

---

# Return Augmented Decision Transformer for Off-Dynamics Reinforcement Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We study offline off-dynamics reinforcement learning (RL) to utilize data from  
2 an easily accessible source domain to enhance policy learning in a target domain  
3 with limited data. Our approach centers on return-conditioned supervised learning  
4 (RCSL), particularly focusing on Decision Transformer (DT) type frameworks,  
5 which can predict actions conditioned on desired return guidance and complete tra-  
6 jectory history. Previous works address the dynamics shift problem by augmenting  
7 the reward in the trajectory from the source domain to match the optimal trajectory  
8 in the target domain. However, this strategy can not be directly applicable in RCSL  
9 owing to (1) the unique form of the RCSL policy class, which explicitly depends  
10 on the return, and (2) the absence of a straightforward representation of the optimal  
11 trajectory distribution. We propose the Return Augmented (REAG) method for DT  
12 type frameworks, where we augment the return in the source domain by aligning  
13 its distribution with that in the target domain. We provide the theoretical analysis  
14 demonstrating that the RCSL policy learned from REAG achieves the same level  
15 of suboptimality as would be obtained without a dynamics shift. We introduce  
16 two practical implementations  $\text{REAG}_{\text{Dara}}^*$  and  $\text{REAG}_{\text{MV}}^*$  respectively. Thorough  
17 experiments on D4RL datasets and various DT-type baselines demonstrate that  
18 our methods consistently enhance the performance of DT type frameworks in  
19 off-dynamics RL.

## 20 1 Introduction

21 Off-dynamics reinforcement learning (Eysenbach et al., 2020; Jiang et al., 2021; Liu et al., 2022; Liu  
22 and Xu, 2024; Guo et al., 2025) arises in decision-making domains such as autonomous driving (Pan  
23 et al., 2017) and medical treatment (Laber et al., 2018; Liu et al., 2023), where direct policy training  
24 through trial-and-error in the target environment is often costly, unethical, or infeasible. A common  
25 strategy is to train the policy in source environments with similar but more accessible dynamics.  
26 However, discrepancies between the source and target environments create a simulation-to-reality  
27 (sim-to-real) gap, which can lead to catastrophic failures when deploying the source-trained policy in  
28 the target environment.

29 Beyond the challenge of dynamics shift, practical scenarios often do not allow real-time online  
30 interaction with the source environment due to time and computational constraints. As a result,  
31 policies must be learned from pre-collected datasets generated by behavior policies. This setting is  
32 particularly difficult, as it combines off-policy, offline, and off-dynamics characteristics. Recently,  
33 supervised learning-based methods (Chen et al., 2021; Brandfonbrener et al., 2022) have emerged as  
34 more stable and scalable alternatives to traditional offline reinforcement learning algorithms grounded  
35 in dynamic programming (Levine et al., 2020). In the offline off-dynamics setting, the majority of  
36 training data is drawn from the source domain, with only a limited portion collected from the target

37 domain. Our study focuses on advancing Decision Transformer (DT) type frameworks (Chen et al.,  
 38 2021; Hu et al., 2024; Zhuang et al., 2024) for off-dynamics reinforcement learning, which can be  
 39 viewed as a special case of return-conditioned supervised learning (RCSL) (Emmons et al., 2021;  
 40 Brandfonbrener et al., 2022). While DT-type methods have gained significant attention across various  
 41 reinforcement learning tasks, no prior work has explicitly tackled the off-dynamics RL problem.

42 There are several previous significant works in off-dynamics reinforcement learning that employ  
 43 reward augmentation to address the dynamics shift between source and target environments (Ey-  
 44 senbach et al., 2020; Liu et al., 2022). In particular, Eysenbach et al. (2020) proposed the DARC  
 45 algorithm to train a policy in the source domain using augmented rewards. These augmentations  
 46 are derived by minimizing the KL distance between the distribution of trajectories generated by the  
 47 learning policy in the source domain and those generated by the optimal policy in the target domain.  
 48 Liu et al. (2022) extended this idea to the offline setting with the DARA algorithm. However, these  
 49 reward augmentation techniques for dynamic programming based RL algorithms are not directly  
 50 applicable to RCSL methods for two primary reasons. First, the policy classes used in RCSL methods  
 51 explicitly depend on the conditional return-to-go function, leading to different trajectory distributions  
 52 that invalidate the trajectory matching methods. Second, the augmentation techniques in Eysenbach  
 53 et al. (2020); Liu et al. (2022) explicitly rely on the form of the optimal trajectory distribution in the  
 54 target domain. In contrast, there is no straightforward representation of the optimal RCSL policy and  
 55 the trajectory distribution. Therefore, novel augmentation mechanisms must be derived for RCSL  
 56 methods to effectively address off-dynamics reinforcement learning.

57 In this work, we propose the Return Augmented (REAG) algorithm, which augments the returns of  
 58 trajectories from the source environment to align with the target environment in DT type framework.  
 59 Through rigorous analysis, we show that the RCSL policy learned with REAG in the source domain  
 60 achieves suboptimality comparable to that learned directly in the target domain without dynamics  
 61 shift. Specifically, our contributions are summarized as follows:

- 62 • We propose a novel method, Return Augmented (REAG), designed specifically for DT-type algo-  
 63 rithms. The approach augments the returns of offline trajectories in the source domain by leveraging  
 64 a small amount of data from the target domain. We develop two practical implementations of REAG:  
 65  $\text{REAG}_{\text{Dara}}^*$ , derived from reward augmentation techniques used in dynamic programming-based  
 66 methods, and  $\text{REAG}_{\text{MV}}^*$  from direct return distribution matching.
- 67 • We provide a rigorous theoretical analysis demonstrating that the return-conditioned policy learned  
 68 from REAG can achieve the same suboptimality as a policy learned directly from the target domain.  
 69 Our analysis relies on the same data coverage assumptions made by Brandfonbrener et al. (2022)  
 70 where there is no dynamics shift, implying that return augmentation could enhance the performance  
 71 of RCSL in off-dynamics RL when the available source dataset size is much larger than the available  
 72 target dataset size.
- 73 • We conduct experiments on the D4RL benchmark by training policies on source datasets collected  
 74 from modified dynamics and evaluating them in the original environments. Across DT-type  
 75 baselines—including DT (Chen et al., 2021), Reinformer (Zhuang et al., 2024) and QT (Hu et al.,  
 76 2024)—both  $\text{REAG}_{\text{Dara}}^*$  and  $\text{REAG}_{\text{MV}}^*$  consistently improve performance, with  $\text{REAG}_{\text{MV}}^*$  showing  
 77 the greatest gains, highlighting the advantage of return-level augmentation.

## 78 2 Preliminary

79 **Sequential Decision-Making.** We consider a general sequential decision-making problem. At  
 80 each step  $t$ , the agent receives an observation  $o_t$  from the environment. Based on the history up  
 81 to step  $t$ , the agent makes action  $a_t$  and receives the reward  $r_t$  from the environment. The agent  
 82 interacts with the environment in episodes with a length  $H$ . We use  $\tau = (o_1, a_1, r_1, \dots, o_H, a_H, r_H)$   
 83 to denote a whole trajectory, and we use  $g(\tau) = \sum_{t=1}^H r_t$  to denote the cumulative return of the  
 84 trajectory. We model the environment as a Markov Decision Process (MDP)  $M$ , which consists of  
 85  $(\mathcal{S}, \mathcal{A}, p, r, H)$ . Here  $\mathcal{S}$  is the state space, each state  $s$  represents the possible history up to some time  
 86 step  $t$ , i.e.,  $s = (o_1, a_1, r_1, \dots, o_t)$ .  $\mathcal{A}$  is the action space,  $p(s'|s, a)$  is the transition dynamics that  
 87 determines the transition probability for the agent to visit state  $s'$  from current state  $s$  with the action  
 88  $a$ .  $r(s, a)$  denotes the reward function. We re-define a trajectory as  $\tau = (s_1, a_1, r_1, \dots, s_H, a_H, r_H)$ .  
 89 We assume that each  $s$  corresponds to one single time step  $t = t(s)$ , and we denote  $g_\pi(s) =$   
 90  $\mathbb{E}_{\tau \sim \pi}[g(\tau)|s_1 = s]$ . Then the goal of the agent is to learn a policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$  that maximizes the  
 91 expected accumulated reward  $J(\pi) := \mathbb{E}_{\tau \sim \pi}[g(\tau)]$ . We denote the optimal policy as  $\pi^*$ .

92 **Offline RL and Decision Transformer.** We consider the offline reinforcement learning setting.  
 93 Given a dataset  $\mathcal{D}$ , the goal of the agent is to learn  $\pi^*$  from  $\mathcal{D}$ . We assume that the trajectories in  $\mathcal{D}$   
 94 are generated from a behavior policy  $\beta$ . In this work, we mainly consider Decision Transformer (DT)  
 95 (Chen et al., 2021) as our backbone algorithm. DT is a type of sequential modeling technique based  
 96 on Transformer (Vaswani et al., 2017) to solve offline RL problems. In detail, DT maintains a function  
 97  $\pi(a|s, g)$  as its policy function. To train the policy  $\pi$ , DT aims to minimize the following negative  
 98 log-likelihood function  $\hat{L}(\pi) := \hat{L}(\pi) := -\sum_{\tau \in \mathcal{D}} \sum_{1 \leq t \leq H} \log \pi(a_t|s_t, g(\tau))$ . To evaluate  $\pi$ , DT  
 99 defines a *conditioning function*  $f : \mathcal{S} \rightarrow \mathbb{R}$ , which maps each state to a return value and guides  
 100 the policy  $\pi_f$  within the environment, where  $\pi_f(a|s) := \pi(a|s, f(s))$ . The conditioning function is  
 101 pivotal in DT, as varying  $f(s)$  for a given state  $s$  results in different policies. To achieve the optimal  
 102 policy,  $f(s)$  should be maximized (Zhuang et al., 2024).

103 **Offline Off-Dynamics RL.** In this work, we consider the offline off-dynamics RL problem, where  
 104 the agent has access to two offline datasets  $\mathcal{D}^S$  and  $\mathcal{D}^T$ .  $\mathcal{D}^S, \mathcal{D}^T$  include the data collected from the  
 105 *source environment*  $M^S$  and the *target environment*  $M^T$ . The source and the target environments  
 106 share the same reward function  $r$ , with different transition dynamics  $p^S$  and  $p^T$ . In practice, we  
 107 assume that the dataset size from the source dataset  $|\mathcal{D}^S|$  is much larger than the data coming from  
 108 the target dataset  $|\mathcal{D}^T|$ . Then the agent aims to find the optimal policy for the target environment  $M^T$   
 109 based on the data from both the source and the target environments. Since the transition dynamics  $p^S$   
 110 and  $p^T$  are different, we can not directly apply existing RL algorithms on the union  $\mathcal{D}^S \cup \mathcal{D}^T$ .

### 111 3 Return Augmentation for Goal Conditioned Supervised Learning

#### 112 3.1 Return-Augmented Framework

113 DT has the potential to address offline off-dynamics reinforcement learning challenges, as shown in  
 114 Table 1. However, it still has certain limitations. To overcome these, we propose a general framework  
 115 that efficiently learns the optimal policy for the target environment using the combined dataset  
 116  $\mathcal{D}^S \cup \mathcal{D}^T$ . Leveraging the return-conditioning nature of DT, we introduce a *return augmentation*  
 117 technique that modifies returns in the offline source dataset through a transformation function. This  
 118 approach allows the policy derived from the augmented source dataset to effectively approximate  
 119 the optimal policy of the target environment, as illustrated in the following equation, where  $\pi^S$   
 120 represents a strong candidate for approximating the optimal policy of the target environment and  $\psi$  is  
 121 the carefully chosen transformation function.

$$\pi^S = \arg \min_{\pi} \hat{L}(\pi) := -\sum_{\tau \in \mathcal{D}^S} \sum_{1 \leq t \leq H} \log \pi(a_t|s_t, \psi(g(\tau))).$$

122 We call our method Return Augmented (REAG) for DT. Next we introduce two methods to construct  
 123  $\psi$ , based on the dynamics-aware reward augmentation (DARA) technique (Eysenbach et al., 2020;  
 124 Liu et al., 2022), and a direct return distribution matching method.

#### 125 3.2 Dynamics-Aware Reward Augmentation

126 We first summarize the idea of DARA. Let  $p^T(s'|s, a)$  denote the transition dynamics of the target  
 127 environment, and  $p^S(s'|s, a)$  denote the source environment. According to the connection of RL and  
 128 probabilistic inference (Levine, 2018), we can turn the optimal policy finding problem into an inference  
 129 problem. We use  $O$  to denote a binary random variable where  $O = 1$  suggests  $\tau$  is a trajectory  
 130 induced by the optimal policy. Given a trajectory  $\tau$ , the likelihood of  $\tau$  being a trajectory induced by  
 131 the optimal policy under the target environment is  $p^T(O = 1|\tau) = \exp(\sum_{t=1}^H r(s_t, a_t)/\eta)$ , where  
 132  $\eta$  is the step size parameter used for tuning. It means that the trajectory with higher cumulative  
 133 rewards is more likely to be the trajectory induced by the optimal policy. We introduce a variational  
 134 distribution  $p_{\pi}^S(\tau) = p(s_1) \prod_{t=1}^T p^S(s_{t+1}|s_t, a_t) \pi(a_t|s_t)$  to approximate  $p_{\pi}^T(O = 1|\tau)$ . Then we  
 135 have

$$\begin{aligned} \log p_{\pi}^T(O = 1) &= \log \mathbb{E}_{\tau \sim p_{\pi}^T(\tau)} p^T(O = 1|\tau) \\ &\geq \mathbb{E}_{\tau \sim p_{\pi}^S(\tau)} [\log p^T(O = 1|\tau) + \log (p_{\pi}^T(\tau)/p_{\pi}^S(\tau))] \\ &= \mathbb{E}_{\tau \sim p_{\pi}^S(\tau)} [\sum_{t=1}^T r(s_t, a_t)/\eta - \log (p^S(s_{t+1}|s_t, a_t)/p^T(s_{t+1}|s_t, a_t))], \end{aligned} \quad (3.1)$$

136 where for the first inequality, we change the distribution of the expectation from  $p_{\pi}^T(\tau)$  to  $p_{\pi}^S(\tau)$  and  
 137 then use Jensen’s inequality to derive the result; the second equation holds due to the assumption/  
 138 modeling that the likelihood of  $\tau$  being a trajectory induced by the optimal policy under the target

139 environment is  $P^T(O = 1|\tau) = \exp(\sum_{t=1}^H r(s_t, a_t)/\eta)$ . Therefore, we obtain an evidence lower  
 140 bound of  $p_\pi^T(O = 1)$ , which equals to find a policy to maximize the value in the source environment,  
 141 with the augmented reward  $r^S(s_t, a_t) = r(s_t, a_t) + \eta \log p^T(s_{t+1}|s_t, a_t) - \eta \log p^S(s_{t+1}|s_t, a_t)$ .  
 142 Following Eysenbach et al. (2020), to estimate the  $\log p^T(s_{t+1}|s_t, a_t) - \log p^S(s_{t+1}|s_t, a_t)$ , we use  
 143 a pair of learned binary classifiers which infers whether the transitions come from the source or  
 144 target environments. Specifically, we denote classifiers  $q_{sas}(\cdot|s, a, s')$  and  $q_{sa}(\cdot|s, a)$ , which return  
 145 the probability for some  $(s, a, s')$  or  $(s, a)$  tuples whether they belong to the source or the target  
 146 environments. Then according to Eysenbach et al. (2020), we have

$$\begin{aligned} \log p^T(s_{t+1}|s_t, a_t) - \log p^S(s_{t+1}|s_t, a_t) &= \Delta r(s_t, a_t, s_{t+1}) \\ &:= \log \frac{q(M^T|s_t, a_t, s_{t+1})}{q(M^S|s_t, a_t, s_{t+1})} - \log \frac{q_{sa}(M^T|s_t, a_t)}{q_{sa}(M^S|s_t, a_t)}. \end{aligned} \quad (3.2)$$

147 For a trajectory  $\tau = (s_1, a_1, r_1, \dots, s_H, a_H, r_H)$ , we denote the transformation  $\psi(g(s_t)) :=$   
 148  $\sum_{h=t}^H r_h + \eta \sum_{h=t}^H \Delta r(s_h, a_h, s_{h+1})$ . We denote such a transformation method as  $\text{REAG}_{\text{Dara}}^*$ .

### 149 3.3 Direct Matching of Return Distributions

150 The reward augmentation strategy in  $\text{REAG}_{\text{Dara}}^*$  stems from the probabilistic inference view of RL  
 151 which matches the distribution of the learning trajectory in the source domain with that of the optimal  
 152 trajectory in the target domain (Eysenbach et al., 2020). However, it does not fully capture the power  
 153 of DT, which is able to induce a *family of policies* that are conditioned on the return-to-go  $f$ . By  
 154 varying  $f$ , DT enables the generation of a diverse range of policies, including the optimal one. In  
 155 contrast,  $\text{REAG}_{\text{Dara}}^*$  assumes a single, fixed target policy, and thus its augmentation strategy cannot  
 156 generalize across multiple policies induced by varying  $f$  in DT. As a result, it cannot find the desired  
 157 return conditioned policy when evaluated with a different  $f$  in the target domain. This motivates us  
 158 to find a return transformation method  $\psi$  to guarantee that  $\pi_f^S(a|s) \approx \pi_f^T(a|s)$  for all  $f$ .

159 We consider a simplified case where both  $D^S$  and  $D^T$  are generated by following the same behavior  
 160 policy  $\beta(a|s)$ . We use  $d_S(A)$  and  $d_T(A)$  to denote the probability for event  $A$  to happen under the  
 161 source and target environments following  $\beta$ . With a slight abuse of notation, we use  $g_S$  and  $g_T$  to  
 162 denote the return following the behavior policy. Then we characterize the learned policies by DT  
 163 under the infinite data regime (Brandfonbrener et al., 2022) for both the source environment and  
 164 target environment. According to Brandfonbrener et al. (2022),  $\pi_f^S(a|s) = P^S(a|s, \psi(g_S) = f(s))$ .  
 165 Then we can express  $\pi^S$  and  $\pi^T$  as

$$\pi_f^S(a|s) = \frac{d_S(a|s)d_S(\psi(g_S) = f(s)|s, a)}{d_S(\psi(g_S) = f(s)|s)}, \quad \pi_f^T(a|s) = \frac{d_T(a|s)d_T(g_T = f(s)|s, a)}{d_T(g_T = f(s)|s)}.$$

166 Since the behavior policies over the source and target environments are the same, we have  $d_S(a|s) =$   
 167  $d_T(a|s)$  for all  $(s, a)$ . Then in order to guarantee  $\pi_f^S(a|s) = \pi_f^T(a|s)$  we only need to guarantee  
 168  $d_S(\psi(g_S(s)) = \cdot|s, a) = d_T(g_T(s) = \cdot|s, a)$ ,  $\forall s, a$ . Denote the cumulative distribution function  
 169 (CDF) of  $g^S$  conditioned on  $s, a$  is  $g^S|s, a \sim G_\beta^S(s, a)$ , and  $g^T|s, a \sim G_\beta^T(s, a)$ . Then if both  
 170  $G_\beta^S(s, a)$  and  $G_\beta^T(s, a)$  are invertible, we can set  $\psi$  as follows

$$\psi(g^S) = G_\beta^{T,-1}(G_\beta^S(g_S; s, a); s, a). \quad (3.3)$$

171 If there exist  $P^S, P^T$ , and  $r$  such that the DARA-type augmented reward-to-go satisfies (3.3), then  
 172 the DARA-type reward augmentation can be deemed as a special case of the transformation (3.3).  
 173 In general,  $G_\beta^T$  and  $G_\beta^S$  are hard to obtain and computationally intractable, making  $\psi$  intractable  
 174 either. We use Laplace approximation to approximate both  $G_\beta^T$  and  $G_\beta^S$  by Gaussian distributions,  
 175 e.g.,  $G_\beta^S(s, a) \sim N(\mu^S(s, a), \sigma_S^2(s, a))$  and  $G_\beta^T(s, a) \sim N(\mu^T(s, a), \sigma_T^2(s, a))$ . We then obtain that

$$\psi(g^S) := \frac{g^S - \mu^S(s, a)}{\sigma^S(s, a)} \cdot \sigma^T(s, a) + \mu^T(s, a). \quad (3.4)$$

176 We denote DT with a  $\psi$  transformation from (3.4) by  $\text{REAG}_{\text{MV}}^*$ , since such a transformation only  
 177 depends on the estimation of mean values  $\mu^S, \mu^T$  and variance  $\sigma^S, \sigma^T$ .

### 178 3.4 Sample Complexity of Off-Dynamics RCSL

179 In this section, we provide an overview of the sample complexity for off-dynamics RCSL. Let  $N^S$   
180 represent the number of trajectories in the source dataset  $\mathcal{D}^S$  and  $N^T$  the number of trajectories in  
181 the target dataset  $\mathcal{D}^T$ . We define  $J^T(\pi)$  as the expected cumulative reward under any policy  $\pi$  within  
182 the target environment. Our theorem is established based on the following assumptions.

183 **Assumption 3.1.** (1) (Return coverage)  $P_\beta^T(g = f(s_1)|s_1) \geq \alpha_f$  for all initial states  $s_1$ . (2) (Near  
184 determinism)  $P(r \neq r(s, a) \text{ or } s' \neq T(s, a)|s, a) \leq \epsilon$  at all  $s, a$  for some functions  $T$  and  $r$ . (3)  
185 (Consistency of  $f$ )  $f(s) = f(s') + r$  for all  $s$ .

186 **Assumption 3.2.** For all  $s$  we assume (1) (Bounded occupancy mismatch)  $P_{\pi_f^{\text{RCSL}}}(s) \leq C_f P_\beta(s)$ ; (2)  
187 (Return coverage)  $P_\beta^T(g = f(s)|s) \geq \alpha_f$ ; and (3) (Domain occupancy overlap)  $d_\beta^T(s) \leq \gamma_f d_\beta^S(s)$ .

188 **Assumption 3.3.** (1) The policy class  $\Pi$  is finite. (2)  $|\log \pi(a|s, g) - \log \pi(a'|s', g')| \leq c$  for  
189 any  $(a, s, g, a', s', g')$  and all  $\pi \in \Pi$ . (3) The approximation error is bounded by  $\epsilon_{\text{approx}}$ , i.e.,  
190  $\min_{\pi \in \Pi} L(\pi) \leq \epsilon_{\text{approx}}$ .

191 **Assumptions 3.1 to 3.3** are the same as the assumptions imposed in Theorem 1, Theorem 2, and  
192 Corollary 3 in [Brandfonbrener et al. \(2022\)](#) respectively. Now we present our theoretical result.

193 **Theorem 3.4.** Under **Assumptions 3.1 to 3.3** on the coverage of the offline dataset and the occupancy  
194 overlap of the source and target environments, with high probability, we have  $J^T(\pi^*) - J^T(\hat{\pi}_f) =$   
195  $O(1/(N^T + N^S)^{1/4})$ , where  $O$  omits terms that are independent of the sample size  $N^T$  of the target  
196 domain and the sample size  $N^S$  of the source domain.

197 **Remark 3.5.** **Theorem 3.4** suggests that the modified samples from the source domain could enhance  
198 the performance of RCSL, for which the sample complexity is approximately  $O((1/N^T)^{1/4})$ .

199 For more theoretical details, please refer to [Appendix B](#).

## 200 4 Experiments

201 In this section, we first outline the fundamental setup of the experiment. We then describe experiments  
202 designed to address specific questions, with each question and its corresponding answer detailed in a  
203 separate subsection.

- 204 • How effective are DT-type methods in mitigating the impact of limited data in target environment?
- 205 • What techniques can be employed to improve the performance of DT-type methods in off-dynamics  
206 scenarios while addressing the constraints of offline data shortages in target environment?
- 207 • How does the performance of DT-type methods compare to baselines in off-dynamics problems?

### 208 4.1 Basic Experiment Setting

209 **Tasks and Environments.** We study established D4RL tasks in the Gym-MuJoCo environment ([Fu](#)  
210 [et al., 2020](#)), a suite built atop the MuJoCo physics simulator, featuring tasks such as locomotion and  
211 manipulation. Particularly, we focused on three environments: Walker2D, Hopper, and HalfCheetah.  
212 In addressing the off-dynamics reinforcement learning problem, we distinguish between the Source  
213 and Target environments. The Target environment is based on the original Gym-MuJoCo framework,  
214 while the Source environment is modified using two shift methods: BodyMass Shift and JointNoise  
215 Shift. In the BodyMass Shift, the mass of the body is altered in the Source environment, whereas in  
216 the JointNoise Shift, random noise is added to the actions.

217 **Dataset.** For the Target Dataset corresponding to the Target Environment, we leverage the official  
218 D4RL data to construct the target datasets: 10T and 1T. The 10T dataset comprises ten times the  
219 number of trajectories compared to the 1T dataset.<sup>1</sup> For the Source Dataset collection, we begin by

<sup>1</sup>Unlike the approach of [Liu et al. \(2022\)](#), which constructs the 1T dataset by selecting the last 1/10 timesteps from the original target dataset (10T), we propose a uniform sampling method across trajectories in the target dataset.

	BEAR			AWR			BCQ			CQL		
	M	M-R	M-E	M	M-R	M-E	M	M-R	M-E	M	M-R	M-E
<b>1T</b>	4.638 $\pm$ 3.882	0.777 $\pm$ 0.105	9.267 $\pm$ 1.692	68.023 $\pm$ 1.687	28.426 $\pm$ 2.974	100.566 $\pm$ 0.513	62.567 $\pm$ 2.459	60.638 $\pm$ 0.683	101.610 $\pm$ 1.309	65.618 $\pm$ 2.818	57.402 $\pm$ 6.161	101.611 $\pm$ 0.143
<b>10T</b>	13.143 $\pm$ 3.016	5.852 $\pm$ 0.168	21.383 $\pm$ 1.237	78.060 $\pm$ 0.772	58.286 $\pm$ 1.684	109.154 $\pm$ 0.976	74.735 $\pm$ 1.184	64.735 $\pm$ 2.555	101.840 $\pm$ 1.962	78.191 $\pm$ 1.839	80.145 $\pm$ 2.286	101.840 $\pm$ 0.467
	MOPO			DT			Reinformer			QT		
	M	M-R	M-E	M	M-R	M-E	M	M-R	M-E	M	M-R	M-E
<b>1T</b>	20.953 $\pm$ 2.715	20.313 $\pm$ 3.488	20.569 $\pm$ 0.983	67.261 $\pm$ 2.316	34.482 $\pm$ 5.890	107.171 $\pm$ 1.611	79.034 $\pm$ 1.506	38.072 $\pm$ 9.174	103.284 $\pm$ 5.437	81.756 $\pm$ 1.671	67.546 $\pm$ 9.505	111.722 $\pm$ 1.398
<b>10T</b>	22.261 $\pm$ 2.811	18.529 $\pm$ 1.760	21.196 $\pm$ 3.103	79.697 $\pm$ 3.348	68.528 $\pm$ 1.924	108.622 $\pm$ 1.815	81.377 $\pm$ 1.903	68.168 $\pm$ 2.661	109.845 $\pm$ 0.726	88.262 $\pm$ 12.886	85.092 $\pm$ 8.727	111.394 $\pm$ 0.469

**Table 1** Performance comparison of algorithms on the **1T**, **10T**, and **1T10S** datasets. In this study, **1T10S(B)** refers to the source dataset under the **BodyMass shift** setting, while **1T10S(J)** corresponds to the source dataset under the **JointNoise shift** setting. Experiments are conducted using the **Medium (M)**, **Medium-Replay (M-R)**, and **Medium-Expert (M-E)** datasets. We present the results for the **Walker2D** environment here; complete results are provided in [appendix E](#). All reported values are averaged over five seeds (0, 1012, 2024, 3036, 4048).

220 modifying the environment through adjustments to the XML file of the MuJoCo simulator. We then  
 221 collect the Random, Medium, Medium-Replay, and Medium-Expert offline datasets in the modified  
 222 environments, following the same data collection procedure as used in D4RL. For further details on  
 223 the dataset collection process and the datasets, please refer to the [Appendix D](#).

224 **Baselines.** In selecting our baseline models, we incorporate a diverse set of well-established off-  
 225 dynamics RL methods, including BEAR (Kumar et al., 2019), AWR (Peng et al., 2019), BCQ  
 226 (Fujimoto et al., 2019), CQL (Kumar et al., 2020), and MOPO (Yu et al., 2020). Furthermore,  
 227 we enhance these baseline models by incorporating DARA augmentation, resulting in augmented  
 228 algorithms that also serve as baselines for comparison with our proposed method. In establishing  
 229 hyperparameters, we ensure consistency across tasks for certain parameters, such as the learning rate  
 230 and the number of iteration steps. Refer to [Appendix D](#) for further details on the parameter settings.

## 231 4.2 Evaluation of Adaptability and Data Efficiency in RCSL Algorithms

232 We evaluate three representative DT-type algorithms include DT (Chen et al., 2021), Reinformer  
 233 (Zhuang et al., 2024) and QT (Hu et al., 2024) to assess their ability to enable an adaptive policy  
 234 while reducing reliance on offline data in the target environment. To conduct this evaluation, we  
 235 perform two experiments: (1) We examine the performance of the three DT-type algorithms under  
 236 varying dataset sizes and quality levels in the target environment; (2) We evaluate their effectiveness  
 237 in off-dynamics scenarios.

238 To assess the impact of dataset size and quality on the performance of DT-type algorithms, we  
 239 evaluate three algorithms using two datasets: a subset of the target data (1T) and the full target  
 240 dataset (10T), comparing the results against other baselines. These experiments aim to quantify  
 241 the performance gap between training on 1T and 10T datasets, highlighting the effects of target  
 242 environment data scarcity and establishing a benchmark for off-dynamics settings. In off-dynamics  
 243 offline RL, instead of relying solely on a large target dataset, we incorporate a small subset of target  
 244 data with a larger source dataset. To examine how effectively algorithms leverage source data, we  
 245 construct the 1T10S dataset by combining a subset of target data (1T) with the full source dataset  
 246 (10S), following the setting of Liu et al. (2022). This dataset serves as the training set for DT-type  
 247 algorithms, whose performance is then evaluated in the target environment. For a comprehensive  
 248 comparison, we benchmark DT-type algorithms against other baseline methods.

249 The evaluation results in [Table 1](#) demonstrate the impact of dataset size and off-dynamics settings on  
 250 algorithm performance. With limited training data, the algorithm’s learning capacity is restricted,  
 251 leading to degraded performance, especially when target-environment data are scarce. To mitigate  
 252 this issue, we incorporate additional source datasets under BodyMass Shift and JointNoise Shift  
 253 settings, which improve generalization to the target environment. However, while leveraging source  
 254 data can partially compensate for the shortage of target data, it remains less effective than training  
 255 with sufficient target-environment data. To further improve DT-type frameworks under off-dynamics  
 256 settings, we propose two return-based augmentation methods,  $\text{REAG}_{\text{MV}}^*$  and  $\text{REAG}_{\text{Dara}}^*$ , which  
 257 can be applied to DT, Reinformer, and QT frameworks. Specifically, applying  $\text{REAG}_{\text{MV}}^*$  yields  
 258  $\text{REAG}_{\text{MV}}^{\text{DT}}$ ,  $\text{REAG}_{\text{MV}}^{\text{Reinf}}$ , and  $\text{REAG}_{\text{MV}}^{\text{QT}}$ , while applying  $\text{REAG}_{\text{Dara}}^*$  produces  $\text{REAG}_{\text{Dara}}^{\text{DT}}$ ,  $\text{REAG}_{\text{Dara}}^{\text{Reinf}}$ ,  
 259 and  $\text{REAG}_{\text{Dara}}^{\text{QT}}$ , demonstrating the promise of these augmentation techniques in enhancing algorithm  
 260 performance under off-dynamics conditions.

		DT			Reinformer			QT		
		IT10S	REAG <sup>DT</sup> <sub>MV</sub>	REAG <sup>DT</sup> <sub>Dara</sub>	IT10S	REAG <sup>Reinf</sup> <sub>MV</sub>	REAG <sup>Reinf</sup> <sub>Dara</sub>	IT10S	REAG <sup>QT</sup> <sub>MV</sub>	REAG <sup>QT</sup> <sub>Dara</sub>
M	BM	78.768±1.233	80.857±1.715↑	78.257±2.423↓	80.857±0.509	82.354±1.479↑	80.666±0.505↓	84.325±0.425	84.582±0.507↑	83.068±0.859↓
	JN	71.068±1.022	75.008±1.834↑	71.779±1.706↑	74.748±1.721	75.008±0.986↑	74.268±1.341↓	80.621±1.143	80.904±1.502↑	78.672±2.201↓
W2D	BM	73.664±1.920	73.708±1.570↑	67.565±0.799↓	67.032±5.767	50.296±14.211↓	66.658±4.303↓	87.292±0.631	87.491±1.226↑	76.169±7.567↓
	JN	58.255±3.181	55.722±2.653↓	62.226±0.383↑	54.801±3.217	47.591±10.244↓	55.438±4.833↑	82.139±1.029	82.363±4.206↑	79.795±4.708↓
M-E	BM	84.430±0.823	88.235±1.886↑	85.328±0.865↑	83.388±0.806	84.897±1.117↑	83.761±0.735↑	93.082±0.348	92.744±0.499↓	94.578±1.383↑
	JN	115.746±1.116	111.060±2.247↓	111.236±0.914↓	117.360±2.550	118.218±1.460↑	117.765±2.499↑	116.149±1.640	118.564±0.697↑	116.115±1.889↓
M	BM	34.057±0.177	39.435±1.239↑	37.787±1.914↑	51.357±3.713	59.085±2.791↑	51.771±5.322↑	49.516±9.798	51.796±9.971↑	62.262±5.348↑
	JN	70.685±0.726	70.356±3.657↓	78.325±2.522↑	70.340±4.633	72.346±5.877↑	70.466±3.728↑	68.656±7.079	73.987±8.080↑	68.709±12.160↑
Hp	BM	64.216±1.504	66.092±0.233↑	60.393±1.086↓	17.534±6.725	20.952±9.794↑	27.238±12.735↑	69.460±13.948	76.287±7.810↑	82.786±11.992↑
	JN	61.870±0.249	77.825±1.638↑	83.525±1.728↑	41.820±15.773	43.985±5.075↑	52.052±10.035↑	93.704±7.559	93.409±4.696↓	51.456±12.168↓
M-E	BM	33.554±0.846	52.873±0.454↑	33.631±1.605↑	68.973±7.512	64.206±12.073↓	73.363±7.674↑	61.162±3.767	73.952±16.294↑	77.279±18.607↑
	JN	108.254±1.583	109.367±1.084↑	108.261±2.612↑	109.256±0.126	109.472±0.103↑	109.255±0.188↓	109.056±0.214	109.803±0.609↑	109.746±0.771↑
M	BM	39.954±0.260	40.250±0.911↑	37.599±0.395↓	37.353±0.483	42.451±0.491↑	38.261±1.238↑	44.656±0.643	47.303±0.318↑	46.383±0.358↑
	JN	47.725±0.431	44.149±3.672↓	47.833±0.284↑	48.274±0.191	43.009±0.307↓	48.404±0.168↑	56.213±0.327	52.394±1.413↓	55.026±0.410↓
Hc	BM	20.966±9.607	27.812±3.256↑	24.059±2.271↑	31.584±1.248	32.114±1.455↑	26.995±4.373	41.300±0.787	42.405±0.729↑	41.359±0.985↑
	JN	36.509±4.414	38.417±4.068↑	38.031±3.529↑	40.296±2.914	40.840±2.880↑	38.436±3.377↓	53.763±0.793	53.870±0.981↑	53.257±0.586↓
M-E	BM	54.981±1.147	56.228±2.930↑	51.357±8.231↓	40.568±0.984	46.048±1.657↑	55.818±1.849↑	71.080±8.802	69.819±5.120↓	76.533±8.022↑
	JN	70.573±8.599	77.762±2.099↑	77.751±2.702↑	76.073±3.878	79.390±0.149↑	78.981±1.198↑	82.961±4.019	83.692±0.699↑	82.148±2.758↓

**Table 2** Performance evaluation of two return augmentation methods,  $\text{REAG}_{\text{MV}}^*$  and  $\text{REAG}_{\text{Dara}}^*$ , integrated with DT, Reinformer, and QT frameworks in off-dynamics scenarios. The experiments are conducted in the Walker2D (W2D), Hopper (Hp), and HalfCheetah (Hc) environments under the Medium (M), Medium-Replay (M-R), and Medium-Expert (M-E) settings. The source environment is modified using two shift conditions: **BodyMass shift (BM)** and **JointNoise shift (JN)**. For reference, the table also includes the performance of the original DT-type methods without augmentation, displayed in gray font. Performance changes due to augmentation are indicated with red upward arrows (↑) for improvements and green downward arrows (↓) for degradations compared to the original DT-type methods. All reported values are averaged over five random seeds (0, 1012, 2024, 3036, 4048).

### 261 4.3 Return Augmentation Methods for Off-Dynamics RL

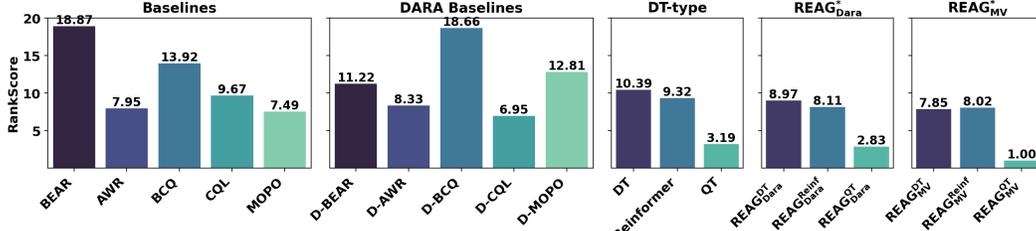
262 Here we discuss how to implement  $\text{REAG}_{\text{MV}}^*$  and  $\text{REAG}_{\text{Dara}}^*$  in practice. We implement  $\text{REAG}_{\text{Dara}}^*$   
263 based on the dynamics-aware reward augmentation method proposed in Liu et al. (2022). For  
264  $\text{REAG}_{\text{MV}}^*$ , it involves training the CQL model across both the Target and Source Environments to  
265 derive the respective value functions, denoted as  $Q_T$  and  $Q_S$ . The derived value functions are then  
266 used to relabel the returns of trajectories in the original dataset. More specifically, the relabeled return  
267  $\hat{g}^S$  is calculated as defined in (3.4). Within this framework, we use  $\mu^S(s, a)$  to denote  $Q_S(s, a)$ ,  
268 and  $Q_T(s, a)$  corresponds to  $\mu^T(s, a)$ . For the computation of  $\sigma_S(s, a)$  and  $\sigma_T(s, a)$ , we employ  
269 the following methodology: For a given state  $s$ , we use the policy of CQL on the source dataset  
270 to obtain  $n$  available actions  $\{a_1^S, a_2^S, \dots, a_n^S\}$  given the state  $s$ , with the corresponding  $Q$  values  
271  $\{Q_S(s, a_1^S), Q_S(s, a_2^S), \dots, Q_S(s, a_n^S)\}$ , and  $n$  available actions  $\{a_1^T, a_2^T, \dots, a_n^T\}$  in the target  
272 environment obtained from the CQL policy trained over the target dataset, with the corresponding  $Q$   
273 values  $\{Q_T(s, a_1^T), Q_T(s, a_2^T), \dots, Q_T(s, a_n^T)\}$ . The standard deviations  $\sigma_S(s, a)$  and  $\sigma_T(s, a)$   
274 are then calculated as specified as follows.

$$\sigma_S(s, a) = \text{std}(Q_S(s, a_1^S), Q_S(s, a_2^S), \dots, Q_S(s, a_n^S)),$$

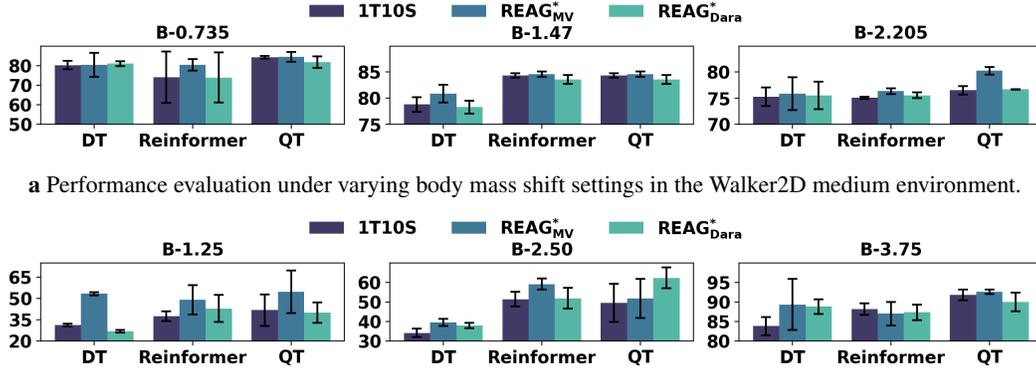
$$\sigma_T(s, a) = \text{std}(Q_T(s, a_1^T), Q_T(s, a_2^T), \dots, Q_T(s, a_n^T)).$$

275 For a detailed discussion, please refer to Section 3. As defined in (3.4), computing the ratio  $\frac{\sigma_T(s, a)}{\sigma_S(s, a)}$   
276 is essential. However, extreme values of this ratio can lead to instability during training. To address  
277 this, we introduce a clipping technique that constrains the ratio within an upper bound  $\theta_1$  and a  
278 lower bound  $\theta_2$ . This helps stabilize  $\text{REAG}_{\text{MV}}^*$  training by mitigating two key challenges. First,  
279 since this ratio depends on the variance of return-to-go in both the source and target environments,  
280 extreme variance values can introduce large gradients or noisy updates, destabilizing training. Second,  
281 variance is estimated using the Q-value function learned through CQL on the source and target  
282 datasets, which may introduce estimation errors. By bounding the ratio within a controlled range,  
283 clipping reduces the impact of these errors and prevents instability.

284 **Table 2** presents a performance comparison of the  $\text{REAG}_{\text{MV}}^*$  and  $\text{REAG}_{\text{Dara}}^*$  return augmentation  
285 techniques integrated into different DT-type frameworks, including DT, Reinformer, and QT, in  
286 off-dynamics scenarios. The results demonstrate that both  $\text{REAG}_{\text{MV}}^*$  and  $\text{REAG}_{\text{Dara}}^*$  effectively  
287 enhance DT-type frameworks, improving performance in most off-dynamics scenarios compared to  
288 their original, non-augmented counterparts. Specifically,  $\text{REAG}_{\text{MV}}^*$ , which augments based on return  
289 values, leverages information from both the source and target environments, making it particularly



**Figure 1** Average normalized rank scores for all baseline algorithms across the Medium, Medium-Replay, and Medium-Expert datasets under BodyMass and JointNoise shift settings in the Walker2D, Hopper, and HalfCheetah environments. Within each setting, algorithms were ranked based on performance, with the top-performing algorithm assigned a rank of 1. Tied scores received the same rank, with subsequent ranks adjusted accordingly. Lower rank scores indicate better overall performance. The original ranks (from 19 algorithms) were normalized to a scale of 1 to 19. The figure presents the average normalized rank scores across the Walker2D, Hopper, and HalfCheetah environments.



a Performance evaluation under varying body mass shift settings in the Walker2D medium environment.

b Performance evaluation under varying body mass shift settings in the Hopper medium environment.

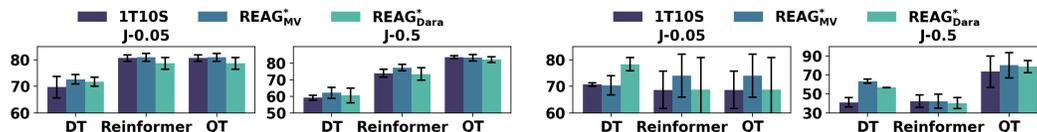
**Figure 2** Performance of  $\text{REAG}_{\text{MV}}^*$  and  $\text{REAG}_{\text{Dara}}^*$  algorithms under different body mass shift settings in the Walker2D and Hopper medium environments. "B-x" denotes that the body mass in the simulator is set to x. The target body mass is 2.94 in the Walker2D environment and 5 in the Hopper environment.

290 well-suited for return-based algorithms. In contrast,  $\text{REAG}_{\text{Dara}}^*$ , which augments based on reward  
 291 values, exhibits more variable performance across different environments and dataset settings. While  
 292  $\text{REAG}_{\text{Dara}}^*$  improves performance in certain cases,  $\text{REAG}_{\text{MV}}^*$  consistently delivers more stable and  
 293 robust improvements.

294 DARA is a widely adopted approach for addressing off-dynamics RL problems by introducing reward  
 295 augmentation to enhance policy adaptation from a source dataset to a target environment while  
 296 minimizing reliance on extensive target data. It seamlessly integrates with traditional offline RL  
 297 frameworks such as CQL and BCQ. In our evaluation, we compare our proposed methods against  
 298 DARA-based approaches, including both traditional RL frameworks and their DARA-augmented  
 299 variants, as well as DT-type frameworks with and without  $\text{REAG}_{\text{MV}}^*$  and  $\text{REAG}_{\text{Dara}}^*$  augmentation,  
 300 providing a comprehensive assessment of augmentation techniques for off-dynamics adaptation. We  
 301 present a comparative ranking where lower average rank scores indicate better overall performance,  
 302 as shown in [Figure 1](#); for the raw results of each setting, please refer to [Appendix E](#). The results  
 303 demonstrate that DT-type frameworks exhibit strong potential in solving off-dynamics RL prob-  
 304 lems, outperforming traditional offline RL methods, particularly in the case of QT. Return-based  
 305 augmentation techniques further enhance effectiveness, with  $\text{REAG}_{\text{MV}}^*$  and  $\text{REAG}_{\text{MV}}^{\text{QT}}$  achieving  
 306 state-of-the-art performance compared to other baselines. Additionally, while DARA effectively  
 307 improves the performance of non-return-based offline RL methods, a noticeable gap remains between  
 308 these approaches and DT-type methods.

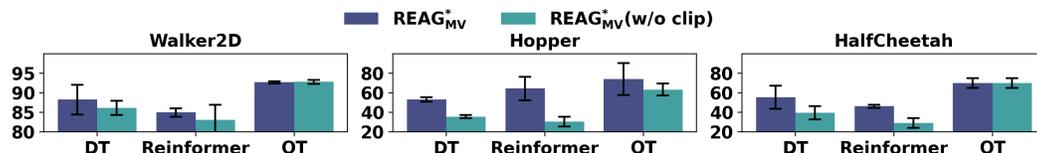
#### 309 4.4 Ablation Studies for Return Augmentation Methods

310 In this section, we present an ablation study to examine the key factors influencing the performance  
 311 of  $\text{REAG}_{\text{Dara}}^*$  and  $\text{REAG}_{\text{MV}}^*$ . We focus on two aspects—**Dynamics Shift** and **Clipped Augmented**



**a** Performance evaluation under varying joint noise shift settings in the Walker2D medium environment. **b** Performance evaluation under varying joint noise shift settings in the Hopper medium environment.

**Figure 3** Performance of  $\text{REAG}_{\text{MV}}^*$  and  $\text{REAG}_{\text{Dara}}^*$  algorithms across varying JointNoise shift settings in the Walker2D and Hopper medium environments. "J-x" denotes the addition of random noise in the range  $(-x, +x)$  to the action.



**Figure 4** Comparison of  $\text{REAG}_{\text{MV}}^*$  with and without the clipping technique in the Medium Expert setting of the Walker2D environment under BodyMass shift. Results are averaged over five seeds.

312 **Return**—while the analysis of **Consistent Augmented Return** and **Return Learning** is deferred to  
 313 [Appendix F](#) due to space limitations.

314 **Dynamics Shift.** To evaluate the impact of shifting source environments on  $\text{REAG}_{\text{MV}}^*$  and  
 315  $\text{REAG}_{\text{Dara}}^*$ , we assess their performance under various BodyMass and JointNoise shift settings.  
 316 The experimental results are presented in [Figure 2](#) and [Figure 3](#). Our findings indicate that as the  
 317 body mass shift increases—creating a greater discrepancy from the target environment—performance  
 318 deteriorates in both the Walker2D and Hopper Medium environments. Similarly, introducing higher  
 319 levels of action noise leads to a decline in performance, suggesting that increased random noise raises  
 320 the likelihood of failure, ultimately resulting in poorer outcomes. This performance degradation is  
 321 particularly evident in the DT framework, highlighting its sensitivity to off-dynamics shifts, whereas  
 322 Reinformer and QT demonstrate greater robustness. Across all shift experiments,  $\text{REAG}_{\text{MV}}^*$  con-  
 323 sistentlly outperforms  $\text{REAG}_{\text{Dara}}^*$ , with the performance gap becoming especially pronounced under  
 324 larger shifts, such as in the Hopper environment with a body mass shift of 1.25.

325 **Clipped Augmented Return.** For data augmentation in  $\text{REAG}_{\text{MV}}^*$ , we apply a clipping technique  
 326 to prevent the occurrence of extreme values. To evaluate its impact, we compare the performance  
 327 of  $\text{REAG}_{\text{MV}}^*$  with and without clipping in the Walker2D, Hopper, and HalfCheetah environments  
 328 under BodyMass shifts with Medium Expert dataset. The results, presented in [Figure 4](#), demonstrate  
 329 that mitigating extreme values generally enhances the performance of  $\text{REAG}_{\text{MV}}^*$ . Additionally, we  
 330 observe that for  $\text{REAG}_{\text{MV}}^{\text{QT}}$ , clipping does not yield significant improvements compared to DT and  
 331 Reinformer. We hypothesize that this is due to the QT mechanism, which inherently regularizes the  
 332 return, whereas DT and Reinformer lack such a mechanism.

## 333 5 Conclusion and Future Work

334 In this work, we introduced the Return-Augmented (REAG) method to improve Decision Trans-  
 335 former-type approaches in off-dynamics reinforcement learning by aligning source-domain returns  
 336 with the target environment. We developed two practical variants,  $\text{REAG}_{\text{Dara}}^*$  and  $\text{REAG}_{\text{MV}}^*$ , and  
 337 provided theoretical guarantees showing that REAG trained on source data can achieve the same  
 338 suboptimality as policies trained directly on target data. Empirical results confirm that REAG en-  
 339 hances DT-type baselines and outperforms several dynamic programming-based methods. Overall,  
 340 REAG offers a promising direction for leveraging source-domain data to address challenges in offline,  
 341 off-policy, and off-dynamics RL. Future work may extend REAG to more diverse environments and  
 342 further refine its augmentation strategies.

## 343 References

- 344 BRANDFONBRENER, D., BIETTI, A., BUCKMAN, J., LAROCHE, R. and BRUNA, J. (2022). When  
345 does return-conditioned supervised learning work for offline reinforcement learning? *Advances in*  
346 *Neural Information Processing Systems* **35** 1542–1553. [1](#), [2](#), [4](#), [5](#), [12](#), [13](#), [14](#)
- 347 CHEN, L., LU, K., RAJESWARAN, A., LEE, K., GROVER, A., LASKIN, M., ABBEEL, P., SRINIVAS,  
348 A. and MORDATCH, I. (2021). Decision transformer: Reinforcement learning via sequence  
349 modeling. *Advances in neural information processing systems* **34** 15084–15097. [1](#), [2](#), [3](#), [6](#), [12](#), [15](#)
- 350 EMMONS, S., EYSENBACH, B., KOSTRIKOV, I. and LEVINE, S. (2021). Rvs: What is essential for  
351 offline rl via supervised learning? *arXiv preprint arXiv:2112.10751* . [2](#), [12](#)
- 352 EYSENBACH, B., ASAWA, S., CHAUDHARI, S., LEVINE, S. and SALAKHUTDINOV, R. (2020).  
353 Off-dynamics reinforcement learning: Training for transfer with domain classifiers. *arXiv preprint*  
354 *arXiv:2006.13916* . [1](#), [2](#), [3](#), [4](#), [12](#)
- 355 FU, J., KUMAR, A., NACHUM, O., TUCKER, G. and LEVINE, S. (2020). D4rl: Datasets for deep  
356 data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219* . [5](#)
- 357 FUJIMOTO, S., MEGER, D. and PRECUP, D. (2019). Off-policy deep reinforcement learning without  
358 exploration. In *International conference on machine learning*. PMLR. [6](#), [15](#)
- 359 GUI, H., PANG, S., YU, S., QIAO, S., QI, Y., HE, X., WANG, M. and ZHAI, X. (2023). Cross-  
360 domain policy adaptation with dynamics alignment. *Neural Networks* **167** 104–117. [12](#)
- 361 GUO, Y., WANG, Y., SHI, Y., XU, P. and LIU, A. (2025). Off-dynamics reinforcement learning via  
362 domain adaptation and reward augmented imitation. *Advances in Neural Information Processing*  
363 *Systems* **37** 136326–136360. [1](#)
- 364 HU, S., FAN, Z., HUANG, C., SHEN, L., ZHANG, Y., WANG, Y. and TAO, D. (2024). Q-value  
365 regularized transformer for offline reinforcement learning. *arXiv preprint arXiv:2405.17098* . [2](#), [6](#),  
366 [15](#)
- 367 JIANG, Y., ZHANG, T., HO, D., BAI, Y., LIU, C. K., LEVINE, S. and TAN, J. (2021). Simgan:  
368 Hybrid simulator identification for domain adaptation via adversarial reinforcement learning. In  
369 *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE. [1](#)
- 370 KUMAR, A., FU, J., SOH, M., TUCKER, G. and LEVINE, S. (2019). Stabilizing off-policy q-learning  
371 via bootstrapping error reduction. *Advances in neural information processing systems* **32**. [6](#), [15](#)
- 372 KUMAR, A., ZHOU, A., TUCKER, G. and LEVINE, S. (2020). Conservative q-learning for offline  
373 reinforcement learning. *Advances in Neural Information Processing Systems* **33** 1179–1191. [6](#), [15](#)
- 374 LABER, E. B., MEYER, N. J., REICH, B. J., PACIFICI, K., COLLAZO, J. A. and DRAKE, J. M.  
375 (2018). Optimal treatment allocations in space and time for on-line control of an emerging  
376 infectious disease. *Journal of the Royal Statistical Society Series C: Applied Statistics* **67** 743–789.  
377 [1](#)
- 378 LEVINE, S. (2018). Reinforcement learning and control as probabilistic inference: Tutorial and  
379 review. *arXiv preprint arXiv:1805.00909* . [3](#)
- 380 LEVINE, S., KUMAR, A., TUCKER, G. and FU, J. (2020). Offline reinforcement learning: Tutorial,  
381 review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643* . [1](#)
- 382 LIU, J., ZHANG, H. and WANG, D. (2022). Dara: Dynamics-aware reward augmentation in offline  
383 reinforcement learning. *arXiv preprint arXiv:2203.06662* . [1](#), [2](#), [3](#), [5](#), [6](#), [7](#), [12](#), [14](#), [15](#)
- 384 LIU, J., ZHANG, Z., WEI, Z., ZHUANG, Z., KANG, Y., GAI, S. and WANG, D. (2024). Beyond  
385 ood state actions: Supported cross-domain offline reinforcement learning. In *Proceedings of the*  
386 *AAAI Conference on Artificial Intelligence*, vol. 38. [12](#)
- 387 LIU, Z., CLIFTON, J., LABER, E. B., DRAKE, J. and FANG, E. X. (2023). Deep spatial q-learning  
388 for infectious disease control. *Journal of Agricultural, Biological and Environmental Statistics* **28**  
389 749–773. [1](#)

- 390 LIU, Z. and XU, P. (2024). Distributionally robust off-dynamics reinforcement learning: Provable  
391 efficiency with linear function approximation. In *International Conference on Artificial Intelligence  
392 and Statistics*. PMLR. 1
- 393 LYU, J., BAI, C., YANG, J.-W., LU, Z. and LI, X. (????). Cross-domain policy adaptation by  
394 capturing representation mismatch. In *Forty-first International Conference on Machine Learning*.  
395 12
- 396 LYU, J., BAI, C., YANG, J.-W., LU, Z. and LI, X. (2024a). Cross-domain policy adaptation by  
397 capturing representation mismatch. In *International Conference on Machine Learning*. PMLR. 12
- 398 LYU, J., XU, K., XU, J., YAN, M., YANG, J., ZHANG, Z., BAI, C., LU, Z. and LI, X. (2024b).  
399 Odr1: A benchmark for off-dynamics reinforcement learning. *arXiv preprint arXiv:2410.20750* .  
400 12
- 401 NIU, H., QIU, Y., LI, M., ZHOU, G., HU, J., ZHAN, X. ET AL. (2022). When to trust your  
402 simulator: Dynamics-aware hybrid offline-and-online reinforcement learning. *Advances in Neural  
403 Information Processing Systems* 35 36599–36612. 12
- 404 PAN, S. J. and YANG, Q. (2009). A survey on transfer learning. *IEEE Transactions on knowledge  
405 and data engineering* 22 1345–1359. 12
- 406 PAN, X., YOU, Y., WANG, Z. and LU, C. (2017). Virtual to real reinforcement learning for  
407 autonomous driving. *arXiv preprint arXiv:1704.03952* . 1
- 408 PENG, X. B., KUMAR, A., ZHANG, G. and LEVINE, S. (2019). Advantage-weighted regression:  
409 Simple and scalable off-policy reinforcement learning. *arXiv preprint arXiv:1910.00177* . 6, 15
- 410 VASWANI, A., SHAZEER, N., PARMAR, N., USZKOREIT, J., JONES, L., GOMEZ, A. N., KAISER,  
411 Ł. and POLOSUKHIN, I. (2017). Attention is all you need. *Advances in neural information  
412 processing systems* 30. 3
- 413 WEN, X., BAI, C., XU, K., YU, X., ZHANG, Y., LI, X. and WANG, Z. (2024). Contrastive  
414 representation for data filtering in cross-domain offline reinforcement learning. In *Forty-first  
415 International Conference on Machine Learning*.  
416 URL <https://openreview.net/forum?id=rReWhol66R> 12
- 417 XU, K., BAI, C., MA, X., WANG, D., ZHAO, B., WANG, Z., LI, X. and LI, W. (2024). Cross-  
418 domain policy adaptation via value-guided data filtering. *Advances in Neural Information Process-  
419 ing Systems* 36. 12
- 420 XU, M., SHEN, Y., ZHANG, S., LU, Y., ZHAO, D., TENENBAUM, J. and GAN, C. (2022). Prompting  
421 decision transformer for few-shot policy generalization. In *international conference on machine  
422 learning*. PMLR. 12
- 423 XUE, Z., CAI, Q., LIU, S., ZHENG, D., JIANG, P., GAI, K. and AN, B. (2024). State regularized  
424 policy optimization on data with dynamics shift. *Advances in neural information processing  
425 systems* 36. 12
- 426 YU, T., THOMAS, G., YU, L., ERMON, S., ZOU, J. Y., LEVINE, S., FINN, C. and MA, T. (2020).  
427 Mopo: Model-based offline policy optimization. *Advances in Neural Information Processing  
428 Systems* 33 14129–14142. 6, 15
- 429 ZHENG, Q., ZHANG, A. and GROVER, A. (2022). Online decision transformer. In *international  
430 conference on machine learning*. PMLR. 12
- 431 ZHUANG, Z., PENG, D., ZHANG, Z., WANG, D. ET AL. (2024). Reinformer: Max-return sequence  
432 modeling for offline rl. *arXiv preprint arXiv:2405.08740* . 2, 3, 6, 15

433 **A Related Work**

434 **Off-dynamics reinforcement learning (RL).** It is a type of domain adaptation problem in RL,  
 435 drawing on concepts from transfer learning (Pan and Yang, 2009). There are many algorithms  
 436 proposed to solve this problem (Niu et al., 2022; Liu et al., 2024; Xu et al., 2024). One of the  
 437 promising approaches is to modify the reward in the source domain. The DARC algorithm (Eysenbach  
 438 et al., 2020) addresses this domain adaptation challenge in the online setting by proposing a reward  
 439 augmentation method that matches the optimal trajectory distribution between the source and target  
 440 domains. Building on this, DARA (Liu et al., 2022) utilizes reward augmentation to supplement a  
 441 limited target dataset with a larger source dataset. Unlike DARC and DARA, which are based on  
 442 dynamic programming, our work adopts the adaptation setting of DARA and introduces a novel  
 443 augmentation method tailored for RCSL, specifically focusing on the Decision Transformer. PAR  
 444 (Lyu et al., 2024a) learns state encoder and state-action encoder utilizing the dynamics representation  
 445 deviation to augment the reward in online settings.

446 **Return Conditioned Supervised Learning (RCSL).** It is a general framework for powerful  
 447 supervised methods in offline RL (Brandfonbrener et al., 2022). Notable works such as RvS (Emmons  
 448 et al., 2021) and Decision Transformer (DT) (Chen et al., 2021) have shown competitive performance  
 449 compared to traditional RL methods. The core idea of RCSL is to condition policies on a desired  
 450 return. In this paper, we primarily focus on DT, which is a specific instance of RCSL and conducts  
 451 offline RL through sequence generation. The generalization potential of DT has inspired researchers  
 452 to explore its use in various settings. For example, Zheng et al. (2022); Xu et al. (2022) leverage  
 453 the DT in the offline-to-online RL and meta RL respectively. However, no prior work has explicitly  
 454 explored the adaptation capabilities of DT in the off-dynamics RL setting.

455 **Additional Related Work.** Niu et al. (2022); Xu et al. (2024); Gui et al. (2023); Lyu et al., 2024b  
 456 present recent advancements in off-dynamics RL methods. Specifically, H2O (Niu et al., 2022)  
 457 performs importance weighting and penalizes Q-values with large dynamics gaps in offline-to-online  
 458 settings. VGDF (Xu et al., 2024) filters data based on value consistency in online off-dynamics RL  
 459 scenarios, while CPD (Gui et al., 2023) employs a dynamics alignment module to minimize discrepan-  
 460 cies. PAR (Lyu et al.) addresses the off-dynamics problem by capturing representation mismatches.  
 461 Lyu et al. (2024b) introduces a newly proposed off-dynamics RL benchmark, demonstrating that IQL  
 462 achieves strong performance in off-dynamics RL settings. For cross-domain offline RL methods,  
 463 BOSA (Liu et al., 2024) tackles OOD state actions with policy optimization and OOD dynamics with  
 464 value optimization, IGDF (Wen et al., 2024) selectively shares transitions from the source domain via  
 465 contrastive learning, and SRPO (Xue et al., 2024) learns the stationary state distribution to regularize  
 466 the policy in a new environment.

467 **B Sample Complexity of Off-Dynamics RCSL**

468 In this section, we provide the rigorous analysis of the sample complexity of the off-dynamics RCSL.  
 469 To this end, we first define some useful notations. We assume there are  $N^S$  trajectories in the source  
 470 dataset  $\mathcal{D}^S$ , and  $N^T$  trajectories in the target dataset  $\mathcal{D}^T$ . Denote  $P_\beta^S$  as the joint distribution of state,  
 471 action, reward and return-to-go induced by the behavior policy  $\beta$  in the source environment, and  $P_\beta^T$   
 472 in the target environment. Denote  $d_\pi^S$  as the marginal distribution of state  $s$  induced by any policy  $\pi$   
 473 in the source environment, and  $d_\pi^T$  in the target environment.

474 Denote  $J^T(\pi)$  as the expected cumulative reward under any policy  $\pi$  and the target environment. For  
 475 any return-to-go  $g$  in the source dataset  $\mathcal{D}^S$ , we transform  $g$  by an oracle defined in (3.3) with others  
 476 remain the same, then we get a modified dataset  $\tilde{\mathcal{D}}^S$ . We denote the mixed dataset as  $\mathcal{D} = \mathcal{D}^T \cup \tilde{\mathcal{D}}^S$ .

477 We first show the sample complexity of DT with only the samples from the target dataset  $\mathcal{D}^T$ . If we  
 478 only use the offline dataset  $\mathcal{D}^T$  collect from the target environment, i.e., at training time we minimizes  
 479 the empirical negative log-likelihood loss:

$$\hat{L}^T(\pi) = - \sum_{\tau \in \mathcal{D}^T} \sum_{1 \leq t \leq H} \log \pi(a_t | s_t, g(s_t)).$$

480 Then we get the following sample complexity guarantee based on the result in Brandfonbrener et al.  
 481 (2022).

482 **Corollary B.1.** There exists a conditioning function  $f : \mathcal{S} \rightarrow \mathbb{R}$  such that assumptions (1)-(3)  
 483 in [Assumption 3.1](#), (1) and (2) in [Assumption 3.2](#) hold. Further assume assumptions (1)-(3) in  
 484 [Assumption 3.3](#) hold. Then for some  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ , we have

$$J^T(\pi^*) - J^T(\hat{\pi}_f) \leq O\left(\frac{C_f}{\alpha_f} H^2 \left(\sqrt{c} \left(\frac{\log |\Pi|/\delta}{N^T}\right)^{1/4} + \sqrt{\epsilon_{\text{approx}}}\right) + \frac{\epsilon}{\alpha_f} H^2\right).$$

485 Now we consider the case of mixed dataset, where we train our policy on both the target dataset and  
 486 the source dataset using the proposed returned conditioned decision transformer methods. Note that  
 487 the size of the target environment dataset is usually small, while the size of the source environment  
 488 dataset is much larger, that is,  $N^T \ll N^S$ . If we incorporate the modified source dataset into the  
 489 supervised learning, that is, we minimize the following empirical negative log-likelihood loss:

$$\hat{L}^{\text{mix}}(\pi) = - \sum_{\tau \in \mathcal{D}} \sum_{1 \leq t \leq H} \log \pi(a_t | s_t, g(s_t)). \quad (\text{B.1})$$

490 An observation is that, with the modified source dataset, the regret  $J^T(\pi^*) - J^T(\hat{\pi}_f)$  can be  
 491 significantly reduced. We state this observation in the following theorem, which is the formal version  
 492 of [Theorem 3.4](#).

493 **Theorem B.2.** There exists a conditioning function  $f$  such that [Assumptions 3.1](#) and [3.2](#) hold. Further  
 494 assume [Assumption 3.3](#) holds. Then for some  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ , we have

$$J^T(\pi^*) - J^T(\hat{\pi}_f) \leq O\left(\frac{C_f}{\alpha_f} \frac{N^S + N^T}{N^S/\gamma_f + N^T} H^2 \left(\sqrt{c} \left(\frac{\log |\Pi|/\delta}{N^T + N^S}\right)^{1/4} + \sqrt{\epsilon_{\text{approx}}}\right) + \frac{\epsilon}{\alpha_f} H^2\right). \quad (\text{B.2})$$

495 **Remark B.3.** Compared to [Corollary B.1](#), [Theorem B.2](#) suggests that the modified samples from  
 496 the source domain could enhance the performance of RCSL when the domain occupancy overlap  
 497 coefficient  $\gamma_f$  is large. In particular, when  $N^S \gg N^T$  and  $\gamma_f = O(1)$ , [\(B.2\)](#) can be simplified to

$$J^T(\pi^*) - J^T(\hat{\pi}_f) \leq O\left(\frac{C_f}{\alpha_f} H^2 \left(\sqrt{c} \left(\frac{\log |\Pi|/\delta}{N^S}\right)^{1/4} + \sqrt{\epsilon_{\text{approx}}}\right) + \frac{\epsilon}{\alpha_f} H^2\right),$$

498 which significantly improves the bound on suboptimality in [Corollary B.1](#).

## 499 C Proof of [Theorem B.2](#)

500 **Lemma C.1** (Corollary 1 of [Brandfonbrener et al. \(2022\)](#)). Under the assumptions in [Assumption 3.1](#),  
 501 there exists a conditioning function such that

$$J^T(\pi^*) - J^T(\pi_f^{\text{RCSL}}) \leq \epsilon \left(\frac{1}{\alpha_f} + 3\right) H^2.$$

502 **Lemma C.2** (Lemma 1 of [Brandfonbrener et al. \(2022\)](#)). For any two policies  $\pi, \pi'$ , we have

$$\|d_\pi^T - d_{\pi'}^T\|_1 \leq 2H \cdot \mathbb{E}_{s \sim d_\pi^T} [TV(\pi(\cdot|s) || \hat{\pi}(\cdot|s))].$$

We define  $d_\beta^{\text{mix}} = \frac{N^T}{N^T + N^S} d_\beta^T + \frac{N^S}{N^T + N^S} d_\beta^S$ . Define

$$L(\hat{\pi}) = \mathbb{E}_{s \sim d_\beta^{\text{mix}}, g \sim P_\beta^T(\cdot|s)} [D_{\text{KL}}(P_\beta^T(\cdot|s, g) || \hat{\pi}(\cdot|s, g))].$$

503 **Theorem C.3.** Consider any function  $f : \mathcal{S} \rightarrow \mathbb{R}$  such that the assumptions in [Assumption 3.2](#) hold.  
 504 Then for any estimated RCSL policy  $\hat{\pi}$  that conditions on  $f$  at test time (denoted by  $\hat{\pi}_f$ ), we have

$$J^T(\pi_f^{\text{RCSL}}) - J^T(\hat{\pi}_f) \leq \frac{C_f \gamma_f}{\alpha_f} H^2 \sqrt{2L(\hat{\pi})}.$$

505 *Proof.* By definition and [Lemma C.2](#), we have

$$\begin{aligned} J^T(\pi_f) - J^T(\hat{\pi}_f) &= H (\mathbb{E}_{P_{\pi_f}^T} [r(s, a)] - \mathbb{E}_{P_{\hat{\pi}_f}^T} [r(s, a)]) \\ &\leq H \cdot \|d_{\pi_f} - d_{\hat{\pi}_f}\|_1 \\ &\leq 2 \cdot \mathbb{E}_{s \sim d_{\pi_f}^T} [TV(\pi_f(\cdot|s) || \hat{\pi}_f(\cdot|s))] H^2. \end{aligned}$$

506 Next, we have

$$\begin{aligned}
& 2 \cdot \mathbb{E}_{s \sim d_{\pi_f}^T} [TV(\pi_f(\cdot|s) || \hat{\pi}_f(\cdot|s))] \\
&= \mathbb{E}_{s \sim d_{\pi_f}^T} \left[ \int_a |P_\beta^T(a|s, f(s)) - \hat{\pi}(a|s, f(s))| \right] \\
&= \mathbb{E}_{s \sim d_{\pi_f}^T} \left[ \frac{P_\beta^T(f(s)|s)}{P_\beta^T(f(s)|s)} \int_a |P_\beta^T(a|s, f(s)) - \hat{\pi}(a|s, f(s))| \right] \\
&\leq 2 \frac{C_f}{\alpha_f} \mathbb{E}_{s \sim d_\beta^T, g \sim P_\beta^T(\cdot|s)} [TV(P_\beta^T(a|s, f(s)) || \hat{\pi}(a|s, f(s)))] \\
&\leq 2 \frac{C_f}{\alpha_f} \frac{N^S + N^T}{N^S/\gamma_f + N^T} \cdot \mathbb{E}_{s \sim d_\beta^{mix}, g \sim P_\beta^T(\cdot|s)} [TV(P_\beta^T(a|s, f(s)) || \hat{\pi}(a|s, f(s)))] \\
&\leq \frac{C_f}{\alpha_f} \frac{N^S + N^T}{N^S/\gamma_f + N^T} \cdot \mathbb{E}_{s \sim d_\beta^{mix}, g \sim P_\beta^T(\cdot|s)} \left[ \sqrt{2KL(P_\beta^T(a|s, f(s)) || \hat{\pi}(a|s, f(s)))} \right] \\
&\leq \frac{C_f}{\alpha_f} \frac{N^S + N^T}{N^S/\gamma_f + N^T} \sqrt{2L(\hat{\pi})}.
\end{aligned}$$

507

□

508 *Proof of Theorem B.2.* Following the same argument in the proof of Corollary 3 in [Brandfonbrener](#)  
509 [et al. \(2022\)](#), we have

$$J^T(\pi_f^{\text{RCSL}}) - J^T(\hat{\pi}_f) \leq O\left(2 \frac{C_f}{\alpha_f} \frac{N^S + N^T}{N^S/\gamma_f + N^T} H^2 \left( \sqrt{c} \left( \frac{\log |\Pi|/\delta}{N^S + N^T} \right)^{1/4} + \sqrt{\epsilon_{\text{approx}}} \right)\right).$$

510 Invoking [Lemma C.1](#), we have

$$J^T(\pi^*) - J^T(\hat{\pi}_f) \leq O\left(2 \frac{C_f}{\alpha_f} \frac{N^S + N^T}{N^S/\gamma_f + N^T} H^2 \left( \sqrt{c} \left( \frac{\log |\Pi|/\delta}{N^T + N^S} \right)^{1/4} + \sqrt{\epsilon_{\text{approx}}} \right) + \frac{\epsilon}{\alpha_f} H^2\right).$$

511 This completes the proof. □

## 512 D Detailed Experiment Setting

### 513 D.1 Environment and Dataset

514 In this section, we provide details of the environments and datasets used in our experiments. We evalu-  
515 ate our approaches in the Hopper, Walker2D, and HalfCheetah environments, using the corresponding  
516 environments from Gym as our target environments.

#### 517 D.1.1 Target Environment Dataset Creation

518 For the target datasets, we construct two distinct datasets: one containing a smaller amount of data  
519 (1T) and another with a larger amount (10T). The 10T dataset consists of ten times the number of  
520 trajectories as the 1T dataset.

521 Both [Liu et al. \(2022\)](#) and our work aim to demonstrate the following two key points:

- 522 • The 10T dataset represents high-quality data, whereas the 1T dataset represents lower-quality  
523 data due to its smaller size.
- 524 • Off-dynamics RL algorithms can enhance performance on 1T by effectively leveraging 10S  
525 source domain data through appropriate data augmentation.

526 [Liu et al. \(2022\)](#) creates the 1T dataset by splitting the original target dataset (10T) based on timesteps,  
527 selecting the last 1/10 timesteps as 1T. However, this approach introduces unintended bias in the  
528 Medium Replay setting, where offline trajectories are collected from a replay buffer in which the  
529 behavior policy improves over time. Consequently, the final 1/10 timesteps tend to exhibit a higher  
530 average return than the overall 10T dataset, undermining the intended quality distinction between 1T  
531 and 10T.

532 To address this issue and ensure a fair evaluation of off-dynamics RL algorithms, we propose a  
 533 uniform sampling method across trajectories in the target dataset. This approach ensures that the  
 534 sampled 1T dataset is a representative subset of the target data, free from biases introduced by  
 535 timestep-based selection. Notably, our method produces a 1T dataset of lower quality than that of Liu  
 536 et al. (2022) in medium replay setting. If an off-dynamics RL algorithm can significantly improve  
 537 performance on our 1T dataset and achieve results comparable to the original 10T dataset, it would  
 538 serve as a more rigorous evaluation and a stronger indicator of the algorithm’s effectiveness.

### 539 D.1.2 Source Environment Dataset Creation

540 We employ BodyMass shift, JointNoise shift to construct the source environments. The following  
 541 descriptions provide detailed insights into the process of creating these source environments.

- 542 • **BodyMass Shift:** The body mass of the agents is modified by adjusting the mass parameters  
 543 in the Gym environment. Detailed body mass settings are provided in Table 3.
- 544 • **JointNoise Shift:** Noise is introduced to the agents’ joints by adding perturbations to the  
 545 actions during source data collection. Specifically, the noise is sampled uniformly from the  
 546 range  $[-0.05, +0.05]$  and applied to the actions when generating the source offline dataset.  
 547 Detailed joint noise settings are provided in Table 3.

548 For the source datasets, we utilize the BodyMass Shift and JointNoise Shift datasets from (Liu et al.,  
 549 2022). Additionally, in our ablation study, we explore variations of BodyMass and JointNoise shifts  
 550 beyond those specified in Table 3. We also collect medium-level source datasets for the Hopper,  
 551 Walker2D, and HalfCheetah environments. Behavior policies are generated by training agents with  
 552 SAC using rlkit (<https://github.com/vitchyr/rlkit>), with checkpoints used for dataset collection. We  
 553 construct the Random, Medium, Medium Replay, and Medium Expert datasets, each reflecting  
 554 different performance levels determined by their corresponding SAC checkpoints. For the JointNoise  
 555 Shift setting, instead of training a new SAC policy and collecting data through environment interaction,  
 we introduce random noise within a specified range directly to the actions.

**Table 3** BodyMass Shift and JointNoise Shift in Hopper, Walker2D and HalfCheetah.

	Hopper		Walker2D		HalfCheetah	
	BodyMass	JointNoise	BodyMass	JointNoise	BodyMass	JointNoise
Source	mass[-1]=2.5	action[-1]+noise	mass[-1]=1.47	action[-1]+noise	mass[4]=0.5	action[-1]+noise
Target	mass[-1]=5.0	action[-1]+0	mass[-1]=2.94	action[-1]+0	mass[4]=1.0	action[-1]+0

556

## 557 D.2 Baselines

558 In our experiments, we use BEAR (Kumar et al., 2019), AWR (Peng et al., 2019), BCQ (Fujimoto  
 559 et al., 2019), CQL (Kumar et al., 2020), and MOPO (Yu et al., 2020), along with their DARA-  
 560 augmented variants (Liu et al., 2022), as baseline methods. We compare these baselines against  
 561 DT (Chen et al., 2021), Reinformer (Zhuang et al., 2024), and QT (Hu et al., 2024), as well as our  
 562 proposed REAG approaches.

## 563 D.3 Hyperparameters

564 In this section, we outline the hyperparameters used for our REAG methodologies. The REAG  
 565 approaches begin with dataset augmentation using either the DARA algorithm ( $REAG_{Dara}^*$ ) or the  
 566 Direct Matching of Return Distributions technique ( $REAG_{MV}^*$ ). The augmented dataset is then  
 567 used to train the DT-type frameworks, which is subsequently evaluated in the target environment.  
 568 Specifically, for  $REAG_{Dara}^*$ , dataset augmentation follows the DARA algorithm, with its corresponding  
 569 hyperparameters provided in Table 4. For  $REAG_{MV}^*$ , the augmentation process is described in  
 570 Section 4.3, where a well-trained Conservative Q-Learning (CQL) model estimates state values,  
 571 incorporating a clipping mechanism to mitigate extreme values. The hyperparameters for CQL  
 572 training are provided in Table 5, the clipping ratios are listed in Table 6, and the training parameters  
 573 for DT, Reinformer, and QT adhere to the settings from their respective original papers.

**Table 4** Hyperparameters used in the DARA algorithm.

Hyperparameter	Value
SA Discriminator MLP Layers	4
SAS Discriminator MLP Layers	4
Hidden Dimension	256
Nonlinearity Function	ReLU
Optimizer	RMSprop
Batch Size	256
Learning Rate	$3 \times 10^{-4}$
$\Delta r$ Coefficient $\eta$	0.1

**Table 5** Hyperparameters used in the CQL algorithm.

Hyperparameter	Value
Actor MLP Layers	3
Critic MLP Layers	3
Hidden Dimension	256
Nonlinearity Function	ReLU
Optimizer	Adam
Batch size	256
Discount Factor	0.99
Temperature	1.0
Actor Learning rate	$1 \times 10^{-4}$
Critic Learning rate	$3 \times 10^{-4}$

**Table 6** Hyperparameters for the Clipping Technique Employed in the REAG<sub>MV</sub>\* Algorithm.

Dataset	Clipping Ratio
Walker2D Random	$0.9 < \theta < 1.25$
Walker2D Medium	$0.9 < \theta < 1.25$
Walker2D Medium Replay	$0.9 < \theta < 1.25$
Walker2D Medium Expert	$0.9 < \theta < 1.25$
Hopper Random	$0.9 < \theta < 1$
Hopper Medium	$0.9 < \theta < 1$
Hopper Medium Replay	$0.9 < \theta < 1$
Hopper Medium Expert	$0.9 < \theta < 1$
HalfCheetah Random	$0.67 < \theta < 1.5$
HalfCheetah Medium	$0.67 < \theta < 1.5$
HalfCheetah Medium Replay	$0.67 < \theta < 1.5$
HalfCheetah Medium Expert	$0.67 < \theta < 1.5$

		BEAR			AWR			BCQ			CQL		
		M	M-R	M-E	M	M-R	M-E	M	M-R	M-E	M	M-R	M-E
W2D	IT	4.638 $\pm$ 3.882	0.777 $\pm$ 0.105	9.267 $\pm$ 1.692	68.023 $\pm$ 1.687	28.426 $\pm$ 2.974	100.566 $\pm$ 0.513	62.567 $\pm$ 2.459	60.638 $\pm$ 0.683	101.610 $\pm$ 1.309	65.618 $\pm$ 2.818	57.402 $\pm$ 6.161	101.611 $\pm$ 0.143
	10T	13.143 $\pm$ 3.016	5.852 $\pm$ 0.168	21.383 $\pm$ 1.237	78.060 $\pm$ 0.772	58.286 $\pm$ 1.684	109.154 $\pm$ 0.976	74.735 $\pm$ 1.184	64.735 $\pm$ 2.555	101.840 $\pm$ 1.962	78.191 $\pm$ 1.839	80.145 $\pm$ 2.286	101.840 $\pm$ 0.467
Hp	IT	8.770 $\pm$ 0.402	5.264 $\pm$ 0.283	31.968 $\pm$ 1.213	55.269 $\pm$ 2.254	54.259 $\pm$ 1.295	54.098 $\pm$ 1.165	63.308 $\pm$ 0.418	68.448 $\pm$ 0.251	62.287 $\pm$ 1.689	74.489 $\pm$ 1.061	71.401 $\pm$ 2.106	82.071 $\pm$ 0.483
	10T	20.398 $\pm$ 2.102	5.554 $\pm$ 0.842	88.236 $\pm$ 2.192	64.494 $\pm$ 2.217	57.548 $\pm$ 1.778	105.361 $\pm$ 1.392	73.462 $\pm$ 2.527	60.385 $\pm$ 0.418	102.775 $\pm$ 1.912	82.945 $\pm$ 0.323	73.168 $\pm$ 2.712	102.071 $\pm$ 1.759
Hc	IT	2.659 $\pm$ 0.167	1.602 $\pm$ 0.275	3.089 $\pm$ 0.104	41.672 $\pm$ 0.732	28.023 $\pm$ 4.027	90.168 $\pm$ 1.398	41.051 $\pm$ 2.908	25.828 $\pm$ 6.142	60.173 $\pm$ 4.175	44.393 $\pm$ 0.263	26.955 $\pm$ 1.274	61.621 $\pm$ 13.093
	10T	10.657 $\pm$ 0.271	19.588 $\pm$ 0.453	16.160 $\pm$ 0.208	42.209 $\pm$ 0.611	41.041 $\pm$ 0.729	90.212 $\pm$ 2.259	46.188 $\pm$ 0.423	38.575 $\pm$ 2.060	95.535 $\pm$ 4.042	49.382 $\pm$ 0.338	46.966 $\pm$ 0.372	87.683 $\pm$ 7.753
		MOPO			DT			Reinformer			QT		
		M	M-R	M-E	M	M-R	M-E	M	M-R	M-E	M	M-R	M-E
W2D	IT	20.953 $\pm$ 2.715	20.313 $\pm$ 3.488	20.569 $\pm$ 0.983	67.261 $\pm$ 2.316	34.482 $\pm$ 5.890	107.171 $\pm$ 1.611	79.034 $\pm$ 1.506	38.072 $\pm$ 9.174	103.284 $\pm$ 5.437	81.756 $\pm$ 1.671	67.546 $\pm$ 9.505	111.722 $\pm$ 1.398
	10T	22.261 $\pm$ 2.811	18.529 $\pm$ 1.760	21.196 $\pm$ 3.103	79.697 $\pm$ 3.348	68.528 $\pm$ 1.924	108.622 $\pm$ 1.815	81.377 $\pm$ 1.903	68.168 $\pm$ 2.661	109.845 $\pm$ 0.726	88.262 $\pm$ 12.886	85.092 $\pm$ 8.727	111.394 $\pm$ 0.469
Hp	IT	31.038 $\pm$ 2.868	5.849 $\pm$ 0.146	35.099 $\pm$ 1.212	66.073 $\pm$ 1.745	61.686 $\pm$ 2.592	100.719 $\pm$ 1.679	74.737 $\pm$ 4.807	36.008 $\pm$ 6.575	60.753 $\pm$ 14.433	70.927 $\pm$ 6.482	83.406 $\pm$ 4.734	108.225 $\pm$ 5.596
	10T	32.769 $\pm$ 1.788	8.638 $\pm$ 1.395	36.161 $\pm$ 2.204	85.589 $\pm$ 5.311	69.701 $\pm$ 5.317	108.087 $\pm$ 1.049	77.792 $\pm$ 4.652	39.856 $\pm$ 12.334	79.389 $\pm$ 28.054	90.176 $\pm$ 0.186	100.321 $\pm$ 1.121	112.908 $\pm$ 3.154
Hc	IT	64.329 $\pm$ 2.096	12.277 $\pm$ 1.953	25.055 $\pm$ 7.834	41.204 $\pm$ 0.430	15.164 $\pm$ 4.847	77.500 $\pm$ 3.323	42.958 $\pm$ 0.065	18.493 $\pm$ 1.584	72.085 $\pm$ 3.491	50.464 $\pm$ 0.127	32.318 $\pm$ 2.435	87.854 $\pm$ 6.657
	10T	65.863 $\pm$ 1.289	59.724 $\pm$ 1.056	28.221 $\pm$ 6.078	42.273 $\pm$ 0.379	34.508 $\pm$ 1.482	82.844 $\pm$ 7.635	43.243 $\pm$ 0.262	39.434 $\pm$ 0.362	87.378 $\pm$ 3.340	51.284 $\pm$ 0.605	49.587 $\pm$ 0.334	94.116 $\pm$ 0.321

**Table 7** Performance comparison of algorithms on the IT and 10T datasets. The experiments are conducted in the Walker2D (W2D), Hopper (Hp), and HalfCheetah (Hc) using the Medium (M), Medium Replay (M-R), and Medium Expert (M-E) datasets. All reported values are averaged over five seeds (0, 1012, 2024, 3036, 4048).

## 574 E Additional Experiments Results

575 This section presents more comprehensive experimental results, including additional variance infor-  
576 mation.

577 In Table 1, we present the partial performance of various algorithms and their DARA variants in  
578 the Walker2D medium environment under BodyMass and JointNoise shift settings, considering  
579 both limited and sufficient target data scenarios. The complete experimental results are provided  
580 in Table 7. Additionally, Table 8 and Table 9 present a comprehensive comparison of different  
581 algorithms and their corresponding augmented variants in addressing the off-dynamics problem  
582 across various environments and shift settings.

## 583 F Ablation Study

584 **Consistent Augmented Return.** It is worth noting that our augmented target returns do not satisfy  
585 the *consistency condition*, which requires that the augmented returns follow  $R_{t+1} - R_t = r_t$ , as  
586 enforced by the original DT. To verify whether consistency is a necessary condition for augmentation  
587 in off-dynamics settings, we conduct the following ablation study. Specifically, we introduce a variant  
588 of REAG<sub>MV</sub>\*, denoted as REAG<sub>MV</sub>\* (consistent), where for each trajectory in the target environment,  
589 return augmentation is applied only to the initial return, while all subsequent augmented returns  
590 are derived using the consistency condition  $R_{t+1} - R_t = r_t$ . The results, presented in Figure 5,  
591 indicate that REAG<sub>MV</sub>\* outperforms its consistency-enforced variant in most cases. This finding  
592 suggests that enforcing consistency does not necessarily improve performance; instead, it may limit  
593 the effectiveness of REAG<sub>MV</sub>\* in the context of off-dynamics offline reinforcement learning.

594 **Return Learning.** To evaluate the learned value functions,  $Q_S$  and  $Q_T$ , and their impact on  
595 REAG<sub>MV</sub>\*, we conduct an ablation experiment. Specifically, we assess the quality of the learned

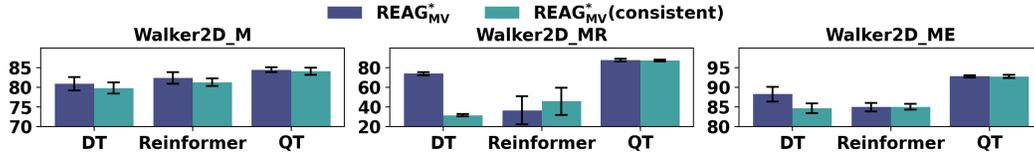
		BEAR	AWR	BCQ	CQL	MOPO	D-BEAR	D-AWR	D-BCQ	D-CQL	D-MOPO
Walker2D M	BM	5.776 ± 1.653	77.442 ± 0.340	70.681 ± 0.539	73.317 ± 1.368	21.617 ± 1.277	6.516 ± 3.220	78.004 ± 0.911	72.023 ± 0.695	74.276 ± 2.582	21.621 ± 1.063
	JN	4.926 ± 1.418	67.636 ± 1.468	62.696 ± 1.037	68.962 ± 0.865	23.552 ± 1.063	6.933 ± 1.884	64.303 ± 0.513	60.681 ± 1.118	69.141 ± 0.944	23.57 ± 0.665
Walker2D M-R	BM	0.0668 ± 4.951	47.033 ± 2.278	50.714 ± 1.918	54.753 ± 0.335	11.563 ± 2.751	1.078 ± 2.083	32.008 ± 1.286	51.447 ± 3.108	57.432 ± 0.764	12.129 ± 2.755
	JN	0.474 ± 0.719	31.623 ± 2.551	50.601 ± 1.611	50.600 ± 1.589	11.379 ± 0.596	0.384 ± 3.823	36.807 ± 2.442	50.714 ± 0.876	51.742 ± 1.061	15.389 ± 0.559
Walker2D M-E	BM	19.799 ± 3.116	110.324 ± 1.053	112.343 ± 1.488	107.187 ± 3.209	18.324 ± 0.708	17.491 ± 2.844	109.743 ± 2.632	113.069 ± 1.602	105.401 ± 2.186	20.741 ± 0.399
	JN	14.225 ± 1.338	104.662 ± 2.370	112.926 ± 1.491	104.019 ± 0.294	17.429 ± 0.639	14.203 ± 1.602	108.915 ± 1.915	111.249 ± 1.092	108.236 ± 1.206	19.325 ± 3.119
Hopper M	BM	22.436 ± 0.103	25.843 ± 0.325	24.853 ± 1.615	49.094 ± 2.207	20.765 ± 3.350	25.608 ± 1.063	26.594 ± 1.267	26.487 ± 1.366	45.101 ± 0.342	21.495 ± 0.848
	JN	8.536 ± 1.965	57.021 ± 0.938	74.559 ± 0.605	71.495 ± 0.126	23.556 ± 1.327	10.576 ± 2.052	61.463 ± 0.702	74.853 ± 0.626	63.611 ± 1.136	24.992 ± 0.944
Hopper M-R	BM	6.282 ± 0.132	55.607 ± 2.310	64.519 ± 0.813	66.455 ± 0.636	5.504 ± 1.701	2.619 ± 0.128	44.883 ± 1.595	64.168 ± 0.291	68.163 ± 0.559	5.482 ± 1.061
	JN	1.841 ± 3.814	37.821 ± 1.205	65.103 ± 0.703	61.302 ± 1.207	5.498 ± 0.568	5.637 ± 0.291	63.937 ± 3.879	64.519 ± 1.102	63.178 ± 1.218	6.147 ± 0.157
Hopper M-E	BM	22.934 ± 3.022	57.595 ± 0.612	109.367 ± 0.834	70.467 ± 2.712	30.541 ± 3.616	31.090 ± 0.463	78.262 ± 0.239	110.014 ± 2.153	72.149 ± 1.934	30.540 ± 0.842
	JN	39.031 ± 1.079	74.708 ± 1.889	108.639 ± 2.028	72.512 ± 0.781	30.537 ± 0.842	33.052 ± 0.385	60.952 ± 0.879	111.587 ± 1.602	94.128 ± 1.213	32.589 ± 1.985
HalfCheetah M	BM	5.431 ± 1.518	42.293 ± 0.862	39.835 ± 0.427	37.081 ± 0.358	58.457 ± 1.449	6.009 ± 1.705	41.800 ± 0.830	39.333 ± 0.506	37.189 ± 0.218	59.311 ± 0.949
	JN	1.948 ± 1.058	41.992 ± 0.762	50.511 ± 0.371	49.046 ± 0.420	61.073 ± 0.315	2.901 ± 0.402	42.545 ± 0.731	52.149 ± 0.457	49.284 ± 0.570	61.447 ± 0.734
HalfCheetah M-R	BM	7.425 ± 1.307	15.988 ± 5.339	32.553 ± 1.258	37.508 ± 0.520	50.429 ± 1.306	4.909 ± 0.562	17.918 ± 3.701	32.095 ± 1.258	37.721 ± 0.440	52.609 ± 0.621
	JN	18.337 ± 0.498	31.742 ± 4.199	46.567 ± 2.563	51.566 ± 0.246	51.918 ± 1.584	17.929 ± 0.479	38.125 ± 1.775	49.066 ± 0.645	52.991 ± 0.438	51.258 ± 1.709
HalfCheetah M-E	BM	4.356 ± 0.431	88.155 ± 1.836	61.771 ± 4.610	61.104 ± 4.131	51.040 ± 4.461	2.948 ± 0.691	89.201 ± 2.419	63.465 ± 3.303	62.665 ± 5.326	56.616 ± 2.609
	JN	3.195 ± 0.391	88.647 ± 2.669	62.486 ± 10.025	84.090 ± 1.109	54.630 ± 10.104	8.789 ± 0.271	89.220 ± 1.800	71.007 ± 4.201	84.210 ± 0.506	60.014 ± 7.011

**Table 8** Performance comparison of traditional offline reinforcement learning algorithms, including BEAR, AWR, BCQ, CQL, and MOPO, along with their DARA-augmented variants, under BodyMass and JointNoise distribution shifts in the Walker2D, Hopper, and HalfCheetah environments. Evaluations are conducted using the Medium (M), Medium Replay (M-R), and Medium Expert (M-E) settings of the 1T10S dataset. The 1T10S dataset comprises a 1T (target) dataset and a 10S (source) dataset. "D-XX" denotes the DARA-augmented variant of the 'XX' algorithm.

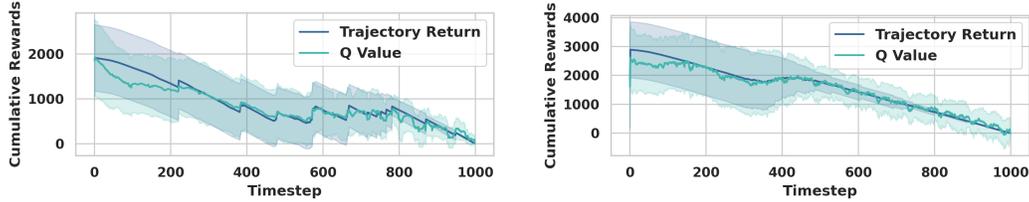
		DT	Reinformer	QT	REAG <sub>MV</sub> <sup>DT</sup>	REAG <sub>MV</sub> <sup>Reinf</sup>	REAG <sub>MV</sub> <sup>QT</sup>	REAG <sub>Dara</sub> <sup>DT</sup>	REAG <sub>Dara</sub> <sup>Reinf</sup>	REAG <sub>Dara</sub> <sup>QT</sup>
Walker2D M	BM	78.768 ± 1.233	80.857 ± 0.509	84.325 ± 0.425	80.857 ± 1.715	82.354 ± 1.479	84.582 ± 0.507	78.257 ± 2.423	80.666 ± 0.505	83.068 ± 0.859
	JN	71.068 ± 1.022	74.748 ± 1.721	80.621 ± 1.143	75.008 ± 1.834	75.008 ± 0.986	80.904 ± 1.502	71.779 ± 1.706	74.268 ± 1.341	78.672 ± 2.201
Walker2D M-R	BM	73.664 ± 1.920	67.032 ± 5.767	87.292 ± 0.631	73.708 ± 1.570	50.296 ± 14.211	87.491 ± 1.226	67.565 ± 0.799	66.658 ± 4.303	76.169 ± 7.567
	JN	58.255 ± 3.181	54.801 ± 3.217	82.139 ± 1.029	55.722 ± 2.653	47.591 ± 10.244	82.363 ± 4.206	62.226 ± 0.383	55.438 ± 4.833	79.795 ± 4.708
Walker2D M-E	BM	84.430 ± 0.823	83.388 ± 0.806	93.082 ± 0.348	88.235 ± 1.886	84.897 ± 1.117	92.744 ± 0.499	85.328 ± 0.865	83.761 ± 0.735	94.578 ± 1.383
	JN	115.746 ± 1.116	117.360 ± 2.550	116.149 ± 1.640	111.060 ± 2.247	118.218 ± 1.460	118.564 ± 0.697	111.236 ± 0.914	117.765 ± 2.499	116.115 ± 1.889
Hopper M	BM	34.057 ± 0.177	51.357 ± 3.713	49.516 ± 9.798	39.435 ± 1.239	59.085 ± 2.791	51.796 ± 9.971	37.787 ± 1.914	51.771 ± 5.322	62.262 ± 5.348
	JN	70.685 ± 0.726	70.340 ± 4.633	68.656 ± 7.079	70.356 ± 3.657	72.346 ± 5.877	73.987 ± 8.080	78.325 ± 2.522	70.466 ± 3.728	68.709 ± 12.160
Hopper M-R	BM	64.216 ± 1.504	17.534 ± 6.725	69.460 ± 13.948	66.092 ± 0.233	20.952 ± 9.794	76.287 ± 7.810	60.393 ± 1.086	27.238 ± 12.735	82.786 ± 11.992
	JN	61.870 ± 0.249	41.820 ± 15.773	93.704 ± 7.559	77.825 ± 1.638	43.985 ± 5.075	93.409 ± 4.696	83.525 ± 1.728	52.052 ± 10.035	51.456 ± 12.168
Hopper M-E	BM	33.554 ± 0.846	68.973 ± 7.512	61.162 ± 3.767	52.873 ± 0.454	64.206 ± 12.073	73.952 ± 16.294	33.631 ± 1.605	73.363 ± 7.674	77.279 ± 18.607
	JN	108.254 ± 1.583	109.256 ± 0.126	109.056 ± 0.214	109.367 ± 1.084	109.472 ± 0.103	109.803 ± 0.609	108.261 ± 2.612	109.255 ± 0.188	109.746 ± 0.771
HalfCheetah M	BM	39.954 ± 0.260	37.353 ± 0.483	44.656 ± 0.643	40.250 ± 0.911	42.451 ± 0.491	47.303 ± 0.318	37.599 ± 0.395	38.261 ± 1.238	46.383 ± 0.358
	JN	47.725 ± 0.431	48.274 ± 0.191	56.213 ± 0.327	44.149 ± 3.672	43.009 ± 0.307	52.394 ± 1.413	47.833 ± 0.284	48.404 ± 0.168	55.026 ± 0.410
HalfCheetah M-R	BM	20.966 ± 9.607	31.584 ± 1.248	41.300 ± 0.787	27.812 ± 3.256	32.114 ± 1.455	42.405 ± 0.729	24.059 ± 2.271	26.995 ± 4.373	41.359 ± 0.985
	JN	36.509 ± 4.414	40.296 ± 2.914	53.763 ± 0.793	38.417 ± 4.068	40.840 ± 2.880	53.870 ± 0.981	38.031 ± 3.529	38.436 ± 3.377	53.257 ± 0.586
HalfCheetah M-E	BM	54.981 ± 1.147	40.568 ± 0.984	71.008 ± 8.802	56.228 ± 2.930	46.048 ± 1.657	69.819 ± 5.120	51.357 ± 8.231	55.818 ± 1.849	76.533 ± 8.022
	JN	70.573 ± 8.599	76.073 ± 3.878	82.961 ± 4.019	77.762 ± 2.099	79.390 ± 0.149	83.692 ± 0.699	77.751 ± 2.702	78.981 ± 1.198	82.148 ± 2.758

**Table 9** Performance comparison of traditional offline reinforcement learning algorithms, including DT, Reinformer and QT, along with our proposed methods REAG<sub>MV</sub><sup>DT</sup>, REAG<sub>Dara</sub><sup>DT</sup>, REAG<sub>MV</sub><sup>Reinf</sup>, REAG<sub>Dara</sub><sup>Reinf</sup>, REAG<sub>MV</sub><sup>QT</sup> and REAG<sub>Dara</sub><sup>QT</sup> under BodyMass and JointNoise distribution shifts in the Walker2D, Hopper, and HalfCheetah environments. Evaluations are conducted using the Medium (M), Medium Replay (M-R), and Medium Expert (M-E) settings of the 1T10S dataset. The 1T10S dataset comprises a 1T (target) dataset and a 10S (source) dataset.

596 value functions in both the source and target domains. We select the Hopper environment with a  
597 medium-expert offline dataset as the target domain and the BodyMass shift as the source domain.  
598 Ideally, the value functions  $Q_S$  and  $Q_T$  learned through REAG<sub>MV</sub><sup>\*</sup> should accurately reflect the returns  
599 of trajectories in their respective domains. To verify this, we train two additional DTs separately on  
600 the source and target offline datasets to obtain policies for these environments. Using these policies,  
601 we generate test trajectories through rollouts and then leverage the learned value functions  $Q_S$  and  
602  $Q_T$ , trained on the 10S and 1T datasets, to predict the returns of these test trajectories. By comparing  
603 the predicted returns with the actual returns, we assess the accuracy of the learned value functions. As  
604 shown in Figure 6, our learned value functions  $Q_S$  and  $Q_T$  accurately reflect the returns of trajectories



**Figure 5** Comparison of REAG<sub>MV</sub> and REAG<sub>MV</sub>(consistent) across Medium, Medium Replay, and Medium Expert settings in the Walker2D environment under BodyMass shift. Results are averaged over five seeds.



**a** Comparison between cumulative rewards and estimated  $Q_S$  values in the source environment with 100 trajectories.

**b** Comparison between cumulative rewards and estimated  $Q_T$  values in the target environment with 100 trajectories.

**Figure 6** Comparison of the cumulative returns and the learned  $Q$ -values for the source (left) and target (right) environments using CQL. Results are plotted with the mean and variance of 100 trajectories.

605 collected by the policies in the source and target environments, demonstrating that the  $Q$ -values used  
 606 in our approach serve as reliable approximations.